

# Unprovability of strong complexity lower bounds in bounded arithmetic

Jiatu Li\*

Institute for Interdisciplinary Information Sciences  
Tsinghua University

Igor C. Oliveira†

Department of Computer Science  
University of Warwick

November 22, 2022

## Abstract

While there has been progress in establishing the unprovability of complexity statements in lower fragments of bounded arithmetic, understanding the limits of Jeřábek’s theory  $\text{APC}_1$  [Jeř07a] and of higher levels of Buss’s hierarchy  $S_2^i$  [Bus86] has been a more elusive task. Even in the more restricted setting of Cook’s theory PV [Coo75], known results often rely on a less natural formalization that encodes a complexity statement using a collection of sentences instead of a single sentence. This is done to reduce the quantifier complexity of the resulting sentences so that standard witnessing results can be invoked.

In this work, we introduce techniques that can establish unprovability results for *stronger theories* and for *sentences of higher quantifier complexity*. In particular, we unconditionally show that  $\text{APC}_1$  cannot prove strong complexity lower bounds separating the third level of the polynomial hierarchy. In more detail, we consider non-uniform average-case separations, and establish that  $\text{APC}_1$  cannot prove a sentence stating that

$$\forall n \geq n_0 \exists f_n \in \Pi_3\text{-SIZE}[n^d] \text{ that is } (1/n)\text{-far from every } \Sigma_3\text{-SIZE}[2^{n^\delta}] \text{ circuit.}$$

This is a consequence of a much more general result showing that, for every  $i \geq 1$ , strong separations for  $\Pi_i\text{-SIZE}[\text{poly}(n)]$  versus  $\Sigma_i\text{-SIZE}[2^{n^{\Omega(1)}}]$  cannot be proved in the theory  $\text{T}_{\text{PV}}^i$  consisting of *all true*  $\forall \Sigma_{i-1}^b$ -sentences in the language of Cook’s theory PV.

A key ingredient in our argument is a convenient *game-theoretic witnessing theorem* that can be applied to sentences of arbitrary quantifier complexity. This theorem is established using methods from proof theory and should be of independent interest. Our unprovability result combines game-theoretic witnessing with extensions of a technique introduced by Krajíček [Kra11] that was recently employed by Pich and Santhanam [PS21] to establish the unprovability of strong complexity lower bounds in PV (i.e., the case  $i = 1$  above, but under a weaker formalization) and in a certain fragment of  $\text{APC}_1$ .

---

\*Email: [lijt19@mails.tsinghua.edu.cn](mailto:lijt19@mails.tsinghua.edu.cn)

†Email: [igor.oliveira@warwick.ac.uk](mailto:igor.oliveira@warwick.ac.uk)

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Results . . . . .	3
1.2	Techniques . . . . .	7
1.3	Organization . . . . .	10
<b>2</b>	<b>Preliminaries</b>	<b>11</b>
<b>3</b>	<b>Auxiliary Results in Logic and Complexity</b>	<b>14</b>
3.1	Total search problems and the polynomial hierarchy . . . . .	14
3.2	The Nisan-Wigderson generator . . . . .	14
3.3	Hardness amplification in the polynomial hierarchy . . . . .	15
3.4	Herbrand's Theorem and the KPT Witnessing Theorem . . . . .	19
3.5	A universal theory for $T_{PV}^i$ . . . . .	20
<b>4</b>	<b>Witnessing Theorems for General Formulas</b>	<b>24</b>
4.1	A game-theoretic witnessing theorem . . . . .	25
4.2	A cut-free sequent calculus . . . . .	27
4.3	Structural transformations of the proof tree . . . . .	29
4.4	Unbounded tree exploration games . . . . .	32
4.5	Partial game trees from sequents . . . . .	34
4.6	Completing the argument: Winning strategies from proofs . . . . .	36
4.7	A special case: Falsifiers with oblivious strategies . . . . .	39
<b>5</b>	<b>Warm-up: Krajíček's Technique and the Pich-Santhanam Result</b>	<b>41</b>
5.1	Formalization of Complexity Lower Bounds . . . . .	42
5.2	Proof of Theorem 5.1 . . . . .	43
5.3	Extensions of the technique and unprovability of weaker lower bounds . . . . .	53
<b>6</b>	<b>Unprovability of Strong Complexity Lower Bounds in Bounded Arithmetic</b>	<b>55</b>
6.1	Unprovability of lower bounds in expressive theories . . . . .	55
6.1.1	Witnessing for $\Pi_i$ vs $\Sigma_i$ lower bounds . . . . .	55
6.1.2	Proof of Theorem 6.1 . . . . .	57
6.1.3	Relaxing the average-case complexity parameter . . . . .	62
6.2	Unprovability of lower bound sentences of higher quantifier complexity . . . . .	63
6.2.1	Witnessing lemma for lower bound sentences . . . . .	63
6.2.2	Proof of Theorem 6.8 . . . . .	64
6.2.3	Relaxing the average-case complexity parameter . . . . .	69
<b>A</b>	<b>Provability in <math>T_{PV}^i</math></b>	<b>74</b>
A.1	Strength of $T_{PV}^i$ and the hierarchy of total functions . . . . .	74
A.2	Strength of $T_{PV}^i$ and the polynomial hierarchy . . . . .	75
A.3	On the provability of $NP \not\subseteq (i.o.)P$ . . . . .	76
<b>B</b>	<b>Model-Theoretic Proof of the KPT Witnessing Theorem for <math>\forall\exists\forall\exists</math> Sentences</b>	<b>77</b>
<b>C</b>	<b>Self-Contained Proof of Theorem 4.20 via Herbrandization</b>	<b>78</b>
<b>D</b>	<b>Lemmas for Hardness Amplification</b>	<b>81</b>
<b>E</b>	<b>The Counting Lemma: Existence of a Good Restriction</b>	<b>85</b>

# 1 Introduction

Establishing unconditional lower bounds on the complexity of computations is one of the primary goals of the theory of computational complexity. While the field has seen progress in the setting of restricted computational devices, such as constant-depth Boolean circuits (e.g., [Hås86, Raz87, Smo87, Wil14]) and monotone Boolean circuits (e.g., [Raz85, And85, AB87]), proving super-linear circuit size lower bounds against general (unrestricted) circuits (see, e.g., [FGHK16, LY22]) and separating complexity classes remain longstanding challenges.

Several barrier results have been proposed to explain why techniques that have been successful in certain settings cannot lead to stronger results. The most well known of them are relativization [BGS75], natural proofs [RR97], and algebrization [AW09] (see also [FLY22, CHO<sup>+</sup>22] for recent examples). While knowledge of these results provides a practical way to check if some approaches are likely to fail, each of these barriers is formulated in an ad-hoc way and is limited in scope. For instance, the natural proofs barrier does not consider a standard notion of “proof” and can be circumvented using simple reductions (see, e.g., [AK10, OS18, CJW19, CHO<sup>+</sup>22]). In general, the aforementioned barriers don’t really tell if we simply haven’t been clever enough to design a better lower bound technique, or if there is a deeper, more fundamental reason for the difficulty of establishing complexity lower bounds and separations.

This motivates the development of a more principled approach to investigate the difficulty of analysing computations and, perhaps more importantly, the intriguing possibility that strong complexity lower bounds might be unprovable from certain mathematical axioms. In order to implement this plan, the first step is to try to understand which logical theories are able to formalise a significant number of results in algorithms and complexity theory. There has been a long and highly successful line of research showing that certain fragments of Peano Arithmetic collectively known as *Bounded Arithmetic* offer a robust class of theories for the formalization of both basic and advanced results in these areas.

*Remark 1.1* (Bounded Arithmetic). Theories of bounded arithmetic aim to capture mathematical proofs that manipulate concepts from a given complexity class (e.g., a proof by induction whose inductive hypothesis can be checked in polynomial time). Notable examples include Cook’s theory PV [Coo75], which formalises polynomial-time reasoning, Jeřábek’s theory APC<sub>1</sub> [Jeř07a], which formalises probabilistic polynomial-time reasoning, and Buss’s theories S<sub>2</sub><sup>i</sup> and T<sub>2</sub><sup>i</sup>, which correspond to the levels of the polynomial-time hierarchy [Bus86].

The correspondence between these theories and the complexity classes is reflected in several ways. For instance, certain *witnessing results* show that every provably total function in a given theory TC (i.e., when  $\forall x \exists! y \varphi(x, y)$  is provable, for some quantifier-free  $\varphi$ ) is computable within the corresponding complexity class  $\mathcal{C}$  (i.e., the function  $y = f(x)$  is in  $\mathcal{C}$ ). There are also close relationships between theories of bounded arithmetic and propositional proof systems, e.g., propositional translations between proofs of certain sentences in PV and S<sub>2</sub><sup>1</sup> and polynomial-size proofs in the extended Frege proof system (see, e.g., [Bey09] and references therein).

Weaker theories corresponding to more fine-grained complexity classes such as TC<sup>0</sup> and NC<sup>1</sup> and the mathematical theorems provable in each of them have also received considerable attention. For instance, key properties of the elementary integer arithmetic operations can be established in theory VTC<sup>0</sup> [Jer22], expander graphs can be constructed and analyzed in theory VNC<sup>1</sup> [BKKK20], and several results from linear algebra can be formalised in theory VNC<sup>2</sup> [TC21]. We refer to Cook and Nguyen [CN10] and Krajíček [Kra95, Kra19] for more information about bounded arithmetic and the logical foundations of complexity theory.

**Complexity Lower Bounds in Bounded Arithmetic.** The study and formalization of complexity lower bounds proofs in bounded arithmetic dates back to Razborov [Raz95b, Raz95a]. We refer to Pich [Pic15a] and to Müller and Pich [MP20] for a comprehensive survey of the area. In particular, the latter paper identifies Jeřábek’s theory APC<sub>1</sub> [Jeř07a] for probabilistic reasoning as a suitable theory for the formalization of several existing circuit lower bounds. (Informally, APC<sub>1</sub> is defined as the extension of Cook’s theory PV

[Coo75] with the dual weak pigeonhole principle for polynomial-time functions.) For instance,  $\text{APC}_1$  can prove super-polynomial lower bounds against bounded-depth circuits and against monotone circuits [MP20], establish the PCP Theorem [Pic15b] (also provable in PV), and formalize randomized matching algorithms [LC11].

Given the expressive power of PV and its extensions, *unconditionally* showing that these theories cannot prove a given result is a non-trivial task. Remarkably, in a recent work, Pich and Santhanam [PS21] employed a technique introduced by Krajíček [Kra11] and further elaborated in [Pic15a] to establish that PV cannot show strong complexity lower bounds separating NP and coNP. More precisely, for each fixed non-deterministic polynomial-time machine  $M$ , they showed that PV cannot prove an average-case lower bound for  $L(M)$  against co-nondeterministic circuits of size  $2^{n^{o(1)}}$ .

In the same work, [PS21] showed that this unprovability result extends to a certain fragment of  $\text{APC}_1$  (see [PS21] for the details and for additional results). They left open the status of the provability of strong complexity lower bounds in  $\text{APC}_1$ . This theory has also been identified in other papers (e.g., [CKKO21]) as an important test case for unconditional unprovability results. This is unsurprising, given the number of advanced results from algorithms and complexity that can be formulated and proved in  $\text{APC}_1$  and its mild extensions (see [Oja04, CN10, Lê14, Pic14, MP20] for many additional examples).

**Witnessing Theorems and Quantifier Complexity.** The approach of [PS21] crucially relies on the KPT Witnessing Theorem [KPT91], a result that can be used to extract computational information from a proof of a sentence with a small number of quantifier alternations. This and similar results (e.g., Herbrand’s Theorem and Buss’s Witnessing Theorem) have been extremely useful tools in unprovability results (see, e.g., [CK07, Kra21, CKKO21]). In order to apply the usual formulation of these witnessing theorems, it is crucial to consider sentences with up to four blocks of quantifiers. In particular, for this reason, the machine  $M$  in the aforementioned result from [PS21] is quantified outside of the sentence (i.e., in the meta-theory). A similar challenge is faced in other papers that consider the unprovability of complexity statements in bounded arithmetic (see, e.g., [KO17] and the subsequent papers [BM20, BKO20]).

**Our Contributions.** We introduce techniques that can establish (unconditional) unprovability results for *stronger theories* and for *sentences of higher quantifier complexity*. We can summarize our main contributions as follows.

- (i) Building on previous works [Kra11, Pic15a, PS21], we establish the unprovability of strong complexity lower bounds in  $\text{APC}_1$  and in more expressive theories of bounded arithmetic. The lower bound sentences showed unprovable refers to separations between the levels of the polynomial hierarchy, where we consider a non-uniform setting and an average-case lower bound against sub-exponential size circuits.
- (ii) We develop a convenient game-theoretic witnessing theorem that allows us to extract computational information from proofs of sentences with an arbitrary number of quantifier alternations. As a consequence, we can consider in Item (i) a more natural formalization of complexity lower bounds compared with [Kra11, Pic15a, PS21]. We believe that this new witnessing result will find more applications in the investigation of the logical foundations of algorithms and complexity theory.

In the next section, we discuss our contributions in detail.

## 1.1 Results

Before formally stating our main unprovability result, we introduce the theories  $T_{PV}^i$  and their common language (vocabulary)  $\mathcal{L}_{PV}$ .

**Theory  $T_{PV}^i$  and Language  $\mathcal{L}_{PV}$ .** We let  $\mathcal{L}_{PV}$  contain the constant symbols 0 and 1, and a function symbol  $f$  for every function in FP, the class of polynomial-time computable functions.<sup>1</sup> In particular,  $\mathcal{L}_{PV}$  contains function symbols for the length function  $|x|$ , addition  $+$ , etc.  $\mathcal{L}_{PV}$  contains the equality predicate  $=$  as its only relation symbol. Note that one can define any polynomial-time computable predicate through its characteristic function, equality, and the constant symbol 1.

For each integer  $i \geq 1$ , we let  $T_{PV}^i$  denote the theory of all true (with respect to the standard model  $\mathbb{N}$ )  $\forall\Sigma_{i-1}^b$  sentences over the language  $\mathcal{L}_{PV}$ .<sup>2</sup> In particular, the theory  $T_{PV}^1$  (which is called  $T_{PV}$  in [PS21]) is at least as strong as Cook's theory PV.<sup>3</sup> We provide some examples of sentences provable in  $T_{PV}^i$  after stating our main result.

**Formalization of Lower Bounds.** In order to consider the provability of a strong complexity lower bound separating the  $i$ -th level of the (non-uniform) polynomial hierarchy, we introduce a sentence  $LB^i(s_1, s_2, m, n_0)$  stating that, for every input length  $n \geq n_0$ , there is a  $\Pi_i$ -circuit  $C$  of size  $\leq s_1(n)$  such that, for every  $\Sigma_i$ -circuit  $D$  of size  $\leq s_2(n)$ , we have

$$\Pr_{x \sim \{0,1\}^n} [C(x) = D(x)] \leq 1 - \frac{m(n)}{2^n}.$$

Here a  $\Pi_i$ -circuit  $C$  (similarly for  $\Sigma_i$  circuits) is simply a standard deterministic Boolean circuit  $C(x, z_1, \dots, z_i)$ , where we define that

$$C(x) = 1 \quad \text{if and only if} \quad \forall z_1 \exists z_2 \dots Q_i z_i C(x, z_1, \dots, z_i) = 1.$$

Formally, let  $\Sigma_i\text{-SIZE}[s(n)]$  and  $\Pi_i\text{-SIZE}[s(n)]$  refer to  $\Sigma_i$ -circuits and  $\Pi_i$ -circuits of size  $s(n)$ , respectively. Let  $LB^i(s_1, s_2, m, n_0)$  denote the following  $\mathcal{L}_{PV}$ -sentence.<sup>4</sup>

$$\begin{aligned} &\forall n \in \text{LogLog with } n \geq n_0 \exists C \in \Pi_i\text{-SIZE}[s_1(n)] \forall D \in \Sigma_i\text{-SIZE}[s_2(n)] \\ &\exists m = m(n) \text{ distinct } n\text{-bit strings } x^1, \dots, x^m \text{ s.t. Error}(C, D, x^i) \text{ for all } i \in [m], \end{aligned}$$

where  $\text{Error}(C, D, x)$  means that the circuits  $C$  and  $D$  do not agree on the input  $x$ . It's easy to see that  $\text{Error}(C, D, x)$  is the disjunction of a  $\Sigma_i^b$ -formula (stating that  $C(x) = 0 \wedge D(x) = 1$ ) and a  $\Pi_i^b$ -formula (stating that  $C(x) = 1 \wedge D(x) = 0$ ). We note that, already for  $i = 1$ ,  $LB^i(s_1, s_2, m, n_0)$  is a  $\forall\Sigma_1^b$ -sentence. In particular, widely used witnessing results such as the KPT Theorem [KPT91] cannot be directly applied to it.

<sup>1</sup>For the reader familiar with bounded arithmetic, we note that in our setup considering polynomial-time functions is equivalent to considering polynomial-time algorithms. See Section 2.2 for more details.

<sup>2</sup>This is a standard class of sentences in bounded arithmetic. Informally, it means that the sentence starts with a block of universal quantifiers, followed by  $i - 1$  blocks of *bounded* quantifiers, i.e.,  $\forall x \leq t$  or  $\exists x \leq t$  for some term  $t$ . The formal definition will be given in Section 2.2.

<sup>3</sup>We use PV to refer to its first-order formalization [Coo75, KPT91], also denoted by  $PV_1$  by some authors.

<sup>4</sup>For the reader that is not familiar with bounded arithmetic, the notation  $n \in \text{LogLog}$  essentially means that all bounded quantifiers refer to objects of length up to  $\text{poly}(2^n)$ . As in [PS21], this makes the unprovability result stronger. Many existing circuit lower bound proofs can be formalized in  $APC_1$  without ever quantifying over objects of length larger than  $\text{poly}(n)$  [MP20].

**Main Unprovability Result.** Next, we state our main theorem on the unprovability of complexity lower bounds in  $T_{PV}^i$  and its corollary for  $APC_1$ .

**Theorem 1.2** (Main Theorem). *For every  $i \geq 1$ ,  $n_0 \in \mathbb{N}$ ,  $\delta \in \mathbb{Q} \cap (0, 1)$ , and  $d \geq 1$ ,*

$$T_{PV}^i \not\vdash LB^i(s_1, s_2, m, n_0),$$

where  $s_1(n) = n^d$ ,  $s_2(n) = 2^{n^\delta}$ , and  $m = 2^n/n$ .

Theorem 1.2 extends the result of [PS21] in two directions. Firstly, it establishes the unprovability of strong complexity lower bounds in theories believed to be much stronger than  $T_{PV}^1$ . Secondly, [PS21] considered a weaker formalization that instead of quantifying over the circuit  $C$  (inside the sentence) considers a collection of sentences  $\{LB_M^1\}_M$ , one for each uniform non-deterministic machine  $M$  (quantified over outside the theory).

*Example 1.3* (The Strength of Theory  $T_{PV}^i$ ). These theories are quite strong already at small values of  $i$ , say  $i = 3$ . Below we give some examples (see Appendix A for a related discussion).

- (i) Fermat's Little Theorem, which states that if  $a^p \not\equiv a \pmod{p}$  then there is  $1 < d < p$  such that  $d \mid p$ , is a true  $\forall \Sigma_1^b$ -sentence in  $\mathcal{L}_{PV}$  and consequently an axiom of  $T_{PV}^2$ . It is unprovable in  $T_{PV}^1$  (therefore also unprovable in PV) unless factoring is easy (see, e.g., [CN10]).
- (ii) The Pigeonhole Principle, which states that for every circuit  $C: [n+1] \rightarrow [n]$  there exists  $x \neq y$  such that  $C(x) = C(y)$ , is also an axiom of  $T_{PV}^2$ . It is not hard to show that even the weaker version of this principle (in which the circuit  $C: [2n] \rightarrow [n]$ ) is unprovable in  $T_{PV}^1$  unless there is no (public-key) collision-resistant hash functions (see, e.g., [Bus08]).
- (iii) The dual Pigeonhole Principle, which states that for every circuit  $C: [n] \rightarrow [n+1]$  there exists  $y \in [n+1]$  such that for all  $x \in [n]$  we have  $C(x) \neq y$ , is in  $T_{PV}^3$ . Even the weak version of this principle (in which the circuit  $C: [n] \rightarrow [2n]$ ) is unprovable in  $T_{PV}^1$  unless EMPTY [Kor21] (also known as Range Avoidance [RSW22]) can be solved in polynomial time with  $O(1)$  circuit-inversion oracle queries.
- (iv) The induction principle for  $\Sigma_i^p$ -predicates is provable in  $T_{PV}^{i+2}$ , while even the induction principle for NP-predicates is unprovable in  $T_{PV}^1$  unless the polynomial-time hierarchy collapses [KPT91, Bus95, Zam96].

Since every axiom of  $APC_1$  is implied by a true  $\forall \Sigma_2^b$ -sentence over the language  $\mathcal{L}_{PV}$  in theory  $T_{PV}^3$  (see Section 2 for the definition of  $APC_1$ ), every sentence provable in  $APC_1$  is also provable in  $T_{PV}^3$ . Consequently, we get the following corollary to Theorem 1.2, which shows that  $APC_1$  cannot establish strong complexity lower bounds separating the third level of the (non-uniform) polynomial hierarchy.

**Corollary 1.4** (Unprovability of Strong Complexity Lower Bounds in  $APC_1$ ). *For every  $n_0 \in \mathbb{N}$ ,  $\delta \in \mathbb{Q} \cap (0, 1)$ , and  $d \geq 1$ ,*

$$APC_1 \not\vdash LB^3(s_1, s_2, m, n_0),$$

where  $s_1(n) = n^d$ ,  $s_2(n) = 2^{n^\delta}$ , and  $m = 2^n/n$ .

Corollary 1.4 establishes the first unconditional result showing the unprovability of strong complexity lower bounds in  $APC_1$ . Previously, [PS21] obtained an extension of their result to a fragment of  $APC_1$ , but left open the provability of the same collection of sentences in  $APC_1$ . Our result is incomparable to theirs in this case, since Corollary 1.4 refers to  $LB^3$  (the third level of the non-uniform polynomial hierarchy) instead of  $\{LB_M^1\}_M$ .



*Remark 1.5* (Relevance to the Logical Foundations of Complexity Theory). The hypothesis that  $P \neq PH$  (which is equivalent to  $P \neq NP$ ) can be interpreted as the statement that polynomial time computations cannot simulate a finite number of bounded quantifier alternations. Our unconditional unprovability result, on the other hand, establishes that  $T_{PV}^i$ , the strongest (sound) theory operating with  $\forall \Sigma_{i-1}^b$  axioms over  $\mathcal{L}_{PV}$ , cannot strongly separate the  $i$ -th level of the polynomial hierarchy.

If the lower bound stated by the  $LB^i$  sentence is true, our result indicates the existence of a fundamental limitation of this theory in reasoning about computations at the  $i$ -th level of the hierarchy and above. In contrast to previous works, which were restricted to subtheories of  $APC_1$ , a significant aspect of Theorem 1.2 is showing that this phenomenon is not caused by a potential weakness of the theory at hand.

**A Game-Theoretic Witnessing Theorem for General Formulas.** In order to prove Theorem 1.2, we introduce a game-theoretic witnessing theorem that can be applied to sentences of high quantifier complexity, such as  $LB^i(s_1, s_2, m, n_0)$ .

For a language (vocabulary)  $\mathcal{L}$ , let  $\varphi(x)$  be a bounded  $\mathcal{L}$ -formula defined as

$$\begin{aligned} \varphi(x) \triangleq & \exists y_1 \leq t_1(x) \forall x_1 \leq s_1(x, y_1) \exists y_2 \leq t_2(x, y_1, x_1) \dots \forall x_{k-1} \leq s_{k-1}(x, y_1, x_1, \dots, y_{k-1}) \\ & \exists y_k \leq t_k(x, y_1, x_1, \dots, y_{k-1}, x_{k-1}) \forall x_k \leq s_k(x, y_1, x_1, \dots, y_k) \phi(x, x_1, \dots, x_k, y_1, \dots, y_k), \end{aligned}$$

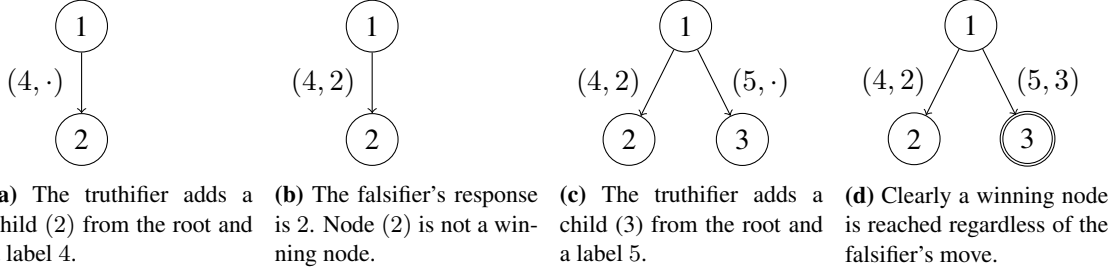
where  $\phi(x, \vec{x}, \vec{y})$  is a quantifier-free  $\mathcal{L}$ -formula. We would like to extract computational information from the provability of  $\forall x \varphi(x)$  in a theory  $\mathcal{T}$ . We achieve this by showing that the provability of this sentence is *equivalent* to the existence of a *winning strategy* in a certain *game*. Moreover, the winning strategy will be computable using *terms* of  $\mathcal{L}$ . Consequently, if the interpretation of each term in a given model  $\mathcal{M}$  of  $\mathcal{T}$  has limited computational complexity, we obtain a *computationally bounded* winning strategy. For simplicity, we discuss the game only informally below, deferring the formal details to Section 4.

We consider an interactive game between two players, the *truthifier* (associated with existential quantifiers in  $\varphi$ ) and the *falsifier* (associated with universal quantifiers in  $\varphi$ ). A *board* is defined as a pair  $(\mathcal{M}, n_0)$ , where  $\mathcal{M}$  is a structure over  $\mathcal{L}$  such that  $\mathcal{M} \models \mathcal{T}$ , and  $n_0 \in \mathcal{M}$  is an element of its domain. The *evaluation game* for the formula  $\varphi(x)$  on the board  $(\mathcal{M}, n_0)$  is played as follows: in the  $i$ -th round of the game ( $1 \leq i \leq k$ ), the truthifier firstly chooses an assignment  $m_i \in \mathcal{M}$  for  $y_i$  such that  $m_i \leq t_i(n_0, m_1, n_1, \dots, m_{i-1}, n_{i-1})$ , then the falsifier chooses an assignment  $n_i \in \mathcal{M}$  for  $x_i$  such that  $n_i \leq s_i(n_0, m_1, n_1, \dots, m_i)$ . The truthifier *wins* if and only if  $\phi(x/n_0, \vec{x}/\vec{n}, \vec{y}/\vec{m})$  holds in  $\mathcal{M}$ . (Note that when playing on a board  $(\mathcal{M}, n_0)$  we set  $x$  in  $\varphi(x)$  to  $n_0$ .)

We will also consider a more general game called the *tree exploration game*. In more detail, we allow the truthifier and falsifier to simultaneously play different evaluation games over the same board  $(\mathcal{M}, n_0)$ . The truthifier has a positional advantage over the falsifier: it can decide where to make the next move, i.e., by either

- (i) making the next move in one of the current games; or
- (ii) starting a new evaluation game over the board  $(\mathcal{M}, n_0)$ ; or
- (iii) playing differently some earlier play, which creates a new game from that position but maintains the existing game plays.

The falsifier must respond to the move of the truthifier in the corresponding evaluation game. Note that the next assignment selected by each player now depends on previous plays in all concurrent games. The truthifier *wins* the tree exploration game if there is a node  $u$  in the current partial game tree that is a winning node for the truthifier, that is, the concatenation of the pairs of elements labelling the edges on the root-to- $u$  path forms a winning transcript of the truthifier in the evaluation game of  $\varphi(x)$  on the board  $(\mathcal{M}, n_0)$ .



**Figure 1:** A transcript of the tree exploration game for  $\varphi(x) = \exists y \leq 2x \forall z < y (y \geq x \wedge (z = 1 \vee z \nmid y))$  (“there is a prime number within  $[x, 2x]$ ”) on the board  $(\mathbb{N}, 3)$ . The truthifier wins by reaching node (3).

The *tree exploration game* of  $\varphi(x)$  is defined as the tree exploration game starting from a partial game tree containing only the root node. See Figure 1 for an example of a transcript of the tree exploration game.

An  $\mathcal{L}$ -strategy of the truthifier in the tree exploration game is described by a sequence of  $\mathcal{L}$ -terms, where each term describes the next move of the truthifier. Finally, a length- $\ell$   $\mathcal{L}$ -strategy is said to be a *universal winning strategy* if the truthifier wins within  $\ell$  moves against all strategies (not necessarily  $\mathcal{L}$ -strategies) of the falsifier on any board  $(\mathcal{M}, n_0)$ . (The “universality” of the strategy comes from the fact that it succeeds over any board  $(\mathcal{M}, n_0)$  and against any strategy of the falsifier. Moreover, the location of the next move of the truthifier in the game tree will be independent of the board and of the strategy of the falsifier.)

Recall that a theory  $\mathcal{T}$  is said to be a universal theory if every axiom of  $\mathcal{T}$  is of the form  $\forall \vec{z} \psi(\vec{z})$ , where  $\psi(\vec{z})$  is a formula free of quantifiers. We show that the provability of the sentence  $\forall x \varphi(x)$  in a universal theory  $\mathcal{T}$  with a certain closure property is equivalent to the existence of a universal winning  $\mathcal{L}$ -strategy of length  $O(1)$  for the truthifier in the tree exploration game of  $\varphi(x)$ .

**Theorem 1.6** (Game-Theoretic Witnessing Theorem). *Let  $\mathcal{T}$  be a universal bounded theory with vocabulary  $\mathcal{L}$  that is closed under if-then-else (see Definition 2.2). Let  $\varphi$  be a bounded  $\mathcal{L}$ -formula of the form*

$$\begin{aligned} \varphi(x) \triangleq & \exists y_1 \leq t_1(x) \forall x_1 \leq s_1(x, y_1) \exists y_2 \leq t_2(x, y_1, x_1) \dots \forall x_{k-1} \leq s_{k-1}(x, y_1, x_1, \dots, y_{k-1}) \\ & \exists y_k \leq t_k(x, y_1, x_1, \dots, y_{k-1}, x_{k-1}) \forall x_k \leq s_k(x, y_1, x_1, \dots, y_k) \phi(x, x_1, \dots, x_k, y_1, \dots, y_k), \end{aligned}$$

where  $\phi(x, \vec{x}, \vec{y})$  is a quantifier-free  $\mathcal{L}$ -formula. Then  $\mathcal{T} \vdash \forall x \varphi(x)$  if and only if there is a universal winning  $\mathcal{L}$ -strategy of length  $O(1)$  for the truthifier in the corresponding tree exploration game of  $\varphi(x)$ .

Beyond its applicability to sentences with an arbitrary number of quantifiers, we stress that two key aspects of Theorem 1.6 are that the winning strategy is computed by  $\mathcal{L}$ -terms and that the truthifier wins in constantly many rounds. (In practice, in order to use this result to obtain computational information from a proof, one typically fixes a particular strategy of the falsifier, which depends on the context and intended application.)

*Remark 1.7.* It is possible to show that Theorem 1.6 is a generalization of the KPT Witnessing Theorem [KPT91]: If the formula  $\varphi(x)$  is an  $\exists\forall$ -formula, the evaluation game for  $\varphi$  has only one round; this means that the tree exploration game for  $\varphi$  is essentially a sequential repetition of the evaluation game (which is equivalent to the Student-Teacher game given by KPT Witnessing Theorem; see Theorem 3.11 and [KPT91, Pic15a]). Indeed, KPT witnessing can also be derived from a less general result that we present in Section 4.7 as a corollary of Theorem 1.6 and that is sufficient for the proof of Theorem 1.2.



## 1.2 Techniques

In this section, we present a high-level overview of the techniques employed to show the game-theoretic witnessing theorem (Theorem 1.6) and the unprovability of strong complexity lower bounds in  $T_{PV}^i$  (Theorem 1.2).

**Game-Theoretic Witnessing Theorem.** Recall that we would like to extract computational information from a proof of a sentence in a bounded theory. Our argument relies on techniques from *proof theory*, a discipline that investigates proofs as precise mathematical objects. In a bit more detail, we explore the structural properties of proofs in Gentzen’s cut-free sequent calculus for (classical) first-order logic. We work with the related proof system G3c defined in the classical textbook [TS00]. Here we only informally discuss the system G3c and its connection to our witnessing theorem, leaving the formal definitions and arguments to Section 4.

Rather than dealing with formulas, the sequent calculus G3c manipulates *sequents*  $\Gamma \Rightarrow \Delta$ , where  $\Gamma$  (called the *antecedent*) and  $\Delta$  (called the *succedent*) are finite multiset of formulas. The semantic interpretation of a sequent  $\Gamma \Rightarrow \Delta$  refers to  $\bigwedge \Gamma \rightarrow \bigvee \Delta$ , i.e., at least one of the formulas in  $\Delta$  holds if all the formulas in  $\Gamma$  hold. The system G3c is known to be *sound* and *complete*, namely a sequent  $\Gamma \Rightarrow \Delta$  is valid (i.e.  $\bigwedge \Gamma \rightarrow \bigvee \Delta$  is true in all first-order structures) if and only if there is a derivation of  $\Gamma \Rightarrow \Delta$  from the rules and axioms of G3c. In particular, assuming that  $\mathcal{T}$  is a universal theory and  $\forall x \varphi(x)$  is a formula provable from  $\mathcal{T}$ , we know by the compactness theorem (see Theorem 4.3) that there is a finite subset  $\Gamma_{\mathcal{T}} \subseteq \mathcal{T}$  such that  $\Gamma_{\mathcal{T}} \vdash \varphi(x)$ , which further means that the sequent  $\Gamma_{\mathcal{T}} \Rightarrow \varphi(x)$  has a G3c-proof. Consequently, we aim to extract computational information from the G3c-proof of  $\Gamma_{\mathcal{T}} \Rightarrow \varphi(x)$ .

The most important feature that makes G3c (and other variants of cut-free sequent calculus) useful in the extraction of information from proofs is that the G3c-proofs are *highly structured*, in the sense that there is essentially no “ad-hoc guess” in G3c when considering intermediate lemmas in a proof of a sentence. Concretely, G3c does not offer rules such as the *Modus Ponens* rule  $\alpha, \alpha \rightarrow \beta \vdash \beta$  that makes the proof of  $\beta$  dependent on a “syntactically irrelevant” guessed formula  $\alpha$ . In fact, each (non-axiom) rule of G3c is uniquely determined by the outer-most connective or quantifier of a chosen occurrence of a formula (which is called the *principal formula* of the rule) in the conclusion. The absence of “ad-hoc guesses” for intermediate formulas in a proof significantly simplifies the analysis of the G3c-proof of  $\Gamma_{\mathcal{T}} \Rightarrow \varphi(x)$  via the syntactic structure of  $\Gamma_{\mathcal{T}}$  and  $\varphi(x)$ . We refer to Example 1.8 for a concrete example.

*Example 1.8.* Consider the propositional rules (L $\rightarrow$ ) and (R $\rightarrow$ ) as an example:

$$\frac{\Gamma \Rightarrow \Delta, \alpha \quad \Gamma, \beta \rightarrow \Delta}{\Gamma, \alpha \rightarrow \beta \Rightarrow \Delta} \text{ (L}\rightarrow\text{)} \quad \frac{\Gamma, \alpha \rightarrow \beta, \Delta}{\Gamma \Rightarrow \alpha \rightarrow \beta, \Delta} \text{ (R}\rightarrow\text{)}$$

The principal formula in each of these rules is the displayed  $\alpha \rightarrow \beta$  in the conclusion. As the name of the rule indicates, (L $\rightarrow$ ) (resp. (R $\rightarrow$ )) is uniquely determined when a formula with  $\rightarrow$  as the outermost connective in the antecedent (resp. succedent) of the conclusion is chosen as the principal formula. In other words, no other rule of the system G3c is admissible once we fix this formula and consider its location in the sequent.

Since  $\mathcal{T}$  is a universal theory, every formula in  $\Gamma_{\mathcal{T}}$  is of the form  $\forall \vec{z} \beta(\vec{z})$  for some quantifier-free formula  $\beta$ . Moreover, by a simple translation we can assume that  $\varphi(x)$  is written in prenex normal form:

$$\varphi(x) = \exists y_1 \forall x_1 \dots \exists y_k \forall x_k \phi(x, x_1, \dots, x_k, y_1, \dots, y_k)$$

for some quantifier-free formula  $\phi$ . Consequently, one can show that the last rule application in the proof of  $\Gamma_{\mathcal{T}} \Rightarrow \phi$  will be uniquely determined by the choice of the principal formula, which is either some

$\forall \vec{z} \beta(\vec{z}) \in \Gamma_{\mathcal{T}}$  or  $\varphi(x)$ . In both cases, we have control over the syntactic structure of the premise of the last rule application. We can then consider the second from the last rule application by looking at all the possible choices of the principal formula, and so on, until we reach the axioms employed in the proof.

By a careful structural induction over the G3c-proof, we demonstrate that the applications of the *quantifier rules* of G3c (omitted in this brief overview) for principal formulas in the succedent have a correspondence to the tree exploration game. Through the applications of these rules and the fact that a (fixed finite) proof establishes the validity of the sentence  $\forall x \varphi(x)$  in all first-order structures that satisfy  $\mathcal{T}$ , we show how to extract a universal winning  $\mathcal{L}$ -strategy of length  $O(1)$  for the truthifier from the proof.

*Example 1.9.* Why does the provability in a universal theory  $\mathcal{T}$  correspond to the tree exploration game instead of the simpler evaluation game? As a conceptual example, one may consider the well-known non-constructive proof of the existence of two irrational numbers  $x, y$  such that  $x^y$  is rational. By the Law of Excluded Middle (i.e.,  $A$  or  $\neg A$ ), one can easily argue that either  $(x, y) = (\sqrt{2}, \sqrt{2})$  or  $(x, y) = ((\sqrt{2})^{\sqrt{2}}, \sqrt{2})$  will be the required pair of irrational numbers. However, we cannot figure out which one of these two possibilities is the correct answer from the structure of this proof. Nevertheless, we can convince any opponent that the original statement is true by a two-round “tree exploration game”: we first propose  $(x, y) = ((\sqrt{2})^{\sqrt{2}}, \sqrt{2})$  and, in case that the opponent argues that  $(\sqrt{2})^{\sqrt{2}}$  is rational, we propose  $(\sqrt{2}, \sqrt{2})$  instead. Similarly, the truthifier’s strategy extracted from the G3c-proof is not guaranteed to witness the existential quantifiers in one shot; it might need to interact with the falsifier for constantly many rounds to produce a correct answer (and each current move of the truthifier can depend on previous moves of both players).

**Unprovability of Strong Complexity Lower Bounds.** We extend the approach of [PS21], which explores a technique from [Kra11, Pic15a]. The main challenge for us is that we must consider the significantly more powerful theory  $T_{PV}^i$  and the (un)provability of a sentence  $LB^i(s_1, s_2, m, n_0)$  with a larger number of quantifier alternations. In particular, while [PS21] considered the provability of a strong complexity lower bound against a fixed machine  $M$ , the sentence  $LB^i(s_1, s_2, m, n_0)$  merely states that there exists a strong separation between  $\Pi_i$  circuits vs  $\Sigma_i$  circuits. This introduces an additional technical difficulty that requires us to also revisit and extend the approach of [Kra11, Pic15a].

Suppose, toward a contradiction, that

$$T_{PV}^i \vdash LB^i(s_1, s_2, m, n_0) ,$$

where  $s_1(n) = n^d$ ,  $s_2(n) = 2^{n^\delta}$ , and  $m = 2^n/n$ . In other words, we assume that the theory  $T_{PV}^i$  proves that for every  $n \geq n_0$  there is a  $\Pi_i$ -circuit  $C_n$  of size  $\leq n^d$  such that, for every  $\Sigma_i$ -circuit  $D_n$  of size  $\leq 2^{n^\delta}$ ,

$$\Pr_{x \sim \{0,1\}^n} [C_n(x) = D_n(x)] \leq 1 - \frac{1}{n}.$$

The key idea behind the argument is that the proof of a strong complexity *lower bound* in bounded arithmetic yields a corresponding complexity *upper bound*. We then argue that the lower bound and the upper bound *contradict each other*. From this, the unprovability of the lower bound sentence follows.

In more detail, our high-level strategy is as follows:

- (i) The provability of the average-case lower bound sentence  $LB^i(s_1, s_2, m, n_0)$  implies the provability in  $T_{PV}^i$  of a *worst-case* lower bound for  $\Pi_i$ -SIZE $[n^d]$  vs  $\Sigma_i$ -SIZE $[2^{n^\delta}]$ . The latter is formalized by a sentence  $LB_{wst}^i(s_1, s_2, n_0)$ .
- (ii) From any  $T_{PV}^i$ -proof of  $LB_{wst}^i(s_1, s_2, n_0)$ , we show how to extract a *complexity upper bound* for an arbitrary  $\Pi_i$ -circuit  $E_m(x)$  over an input  $x$  of length  $m$  and of size at most  $\text{poly}(m)$ . (This is done

outside the theory  $T_{PV}^i$ .) More precisely, we show that there is a deterministic circuit  $B_m$  with  $\Sigma_{i-1}^p$ -oracle gates and of size  $\leq 2^{m^{o(1)}}$  such that

$$\Pr_{x \sim \{0,1\}^m} [E_m(x) = B_m(x)] \geq 1/2 + 2^{-m^{o(1)}}.$$

- (iii) We invoke a hardness amplification result for the (non-uniform) polynomial hierarchy to conclude that, on any large enough input length  $n$ , every  $\Pi_i$ -circuit  $C_n$  of size  $\leq n^d$  agrees with some  $\Sigma_i$ -circuit  $D_n$  of size  $\leq 2^{n^\delta}$  on more than a  $1 - 1/n$  fraction of the inputs. (If this is not the case, we would be able to use hardness amplification to contradict the previous item.)

Since  $T_{PV}^i$  is a *sound* theory, i.e., every theorem of  $T_{PV}^i$  is a true sentence, Item (iii) is in contradiction with the complexity lower bound stated in  $LB^i(s_1, s_2, m, n_0)$ . Consequently,  $T_{PV}^i$  does not prove this sentence.

Item (i) is trivial, since the provability of an average-case lower bound immediately yields the provability of a worst-case lower bound against circuits of the same size. Item (iii) requires an extension of a hardness amplification result of Healy, Vadhan, and Viola [HVV06] to higher levels of the polynomial hierarchy. We verify that this is possible in Section 3.3. The most challenging step of the proof is Item (ii), which we discuss next.

*General upper bounds from the provability of a complexity lower bound.* In Item (ii) we aim to extract computational information from a proof of  $LB_{wst}^i(s_1, s_2, n_0)$  in  $T_{PV}^i$ . For this, we would like to invoke our game-theoretic witnessing theorem (Theorem 1.6). Since this result can only be applied to a *universal theory*, the first step is to introduce a convenient universal theory that is at least as powerful as  $T_{PV}^i$ . Using standard techniques from logic, we construct a universal theory  $UT_{PV}^i$  with all the necessary properties (see Theorem 3.22 in Section 3.5). While the axioms of  $UT_{PV}^i$  are structurally simpler (i.e., universal sentences), the terms of  $UT_{PV}^i$  no longer correspond to polynomial-time functions. However, a careful construction of  $UT_{PV}^i$  ensures that its terms (when interpreted over the standard model) correspond to functions in  $FP_{i-1}^{\Sigma^p}$ , which will be sufficient for our purposes. In addition to the (syntactic) simplification of the axioms of  $T_{PV}^i$ , a benefit of  $UT_{PV}^i$  is that the worst-case lower bound sentence  $LB_{wst}^i(s_1, s_2, n_0)$ , whose quantifier complexity grows with  $i$ , simplifies to a  $\forall \Sigma_4^b$ -sentence  $ULB_{wst}^i(s_1, s_2, n_0)$  in the vocabulary of  $UT_{PV}^i$ . (This quantifier complexity is still too high for the KPT witnessing theorem. We discuss their witnessing theorem and the difficulty of extending it to more quantifiers in Section 3.4 and Appendix B.)

Since  $ULB_{wst}^i(s_1, s_2, n_0)$  is also provable in the universal theory  $UT_{PV}^i$ , we can invoke the game-theoretic witnessing theorem with  $\mathcal{T} = UT_{PV}^i$  and on the formula  $\varphi(x)$  corresponding to  $ULB_{wst}^i(s_1, s_2, n_0)$ . (For this overview, think of  $x$  as the input length  $n$ .) Consequently, there is a universal winning  $\mathcal{L}(UT_{PV}^i)$ -strategy for the truthifier (existential player) in the *tree exploration game* of  $\varphi(x)$ . In particular, for every input length  $n \geq n_0$ , the truthifier has a winning strategy computed by functions in  $FP_{i-1}^{\Sigma^p}$  that succeeds within  $O(1)$  plays in producing a  $\Pi_i$ -circuit  $C_n$  of size  $\leq n^d$  that cannot be computed (in the worst case) by  $\Sigma_i$ -circuits of size  $\leq 2^{n^\delta}$ .

The plan for the remainder of the proof is to fix a *particular strategy of the falsifier*, which will depend on the circuit  $E_m$  from Item (ii) that we would like to approximate, and to show that using the  $FP_{i-1}^{\Sigma^p}$ -computable winning strategy of the truthifier we can obtain a good circuit  $B_m$  for  $E_m$ .

Similarly, in the simpler context of the Student-Teacher game obtained from the KPT Witnessing Theorem and for a worst-case lower bound sentence that refers to a fixed machine  $M$ , [Kra11, Pic15a, PS21] showed that an average-case complexity upper bound follows from the provability of a worst-case lower bound. We provide a simple example of how this can be done in Section 5, when we discuss Student-Teacher games with a single round in the context of [PS21]. For games with more than one round, techniques from

pseudorandomness and a more elaborated strategy that employs the Nisan-Wigderson generator [NW94] play an important role in the argument from [Kra11, Pic15a, PS21].

In our context, the following difficulties arise:

- (1) We need to consider the considerably more complicated tree exploration game played between the truthifier and the falsifier.
- (2) The machine  $M$  becomes an arbitrary circuit  $C'$  that the falsifier proposes as a candidate hard function, and different circuits can be proposed until the winning strategy of the truthifier succeeds in producing a hard circuit  $C_n$ .

We are able to avoid a difficult analysis in Item (1) by considering a simpler setting of the tree exploration game that is sufficient for our purposes. In more detail, when considering the strategy for the falsifier based on the circuit  $E_m$  that we would like to approximate, the play of the falsifier in the current node of the game tree only depends on the partial play of the *evaluation game* corresponding to the moves of both players in the root-to-node path of the *tree exploration game*. We develop this simpler framework in Section 4.7.

Finally, in order to address Item (2), we show that it is possible to modify the use of the Nisan-Wigderson generator in [Kra11, Pic15a] when defining the strategy of the falsifier so that even if the truthifier changes the candidate hard circuit  $O(1)$  times when we execute its winning strategy, we are still able to obtain a non-trivial complexity upper bound for  $E_m$ . We refer to Section 6 for the technical details.

### 1.3 Organization

The remaining sections of the paper are organised as follows:

- Section 2 presents some basic definitions in logic and complexity and fixes notation.
- Section 3 establishes several auxiliary results needed for the proof of Theorem 1.2.
- Section 4 states and proves the most general form of our game-theoretic witnessing result (Theorem 1.6). A simpler version that is sufficient for the proof of Theorem 1.2 is derived in Section 4.7.
- Section 5 provides an exposition of Krajíček’s technique [Kra11] (further elaborated in [Pic15a]) and of the main unprovability result from Pich and Santhanam [PS21] in a language that will be more convenient when discussing our proofs.
- Section 6 combines and extends results from the previous sections in order to establish Theorem 1.2.
- Appendix A discusses provability in the theories  $T_{PV}^i$  and relates their strength to certain computational assumptions.
- Appendix B presents a standard model-theoretic proof of the KPT Witnessing Theorem [KPT91].
- Appendix C provides a self-contained proof of the witnessing theorem presented in Section 4.7 using Herbrandization instead of sequent calculus.
- Appendix D and Appendix E prove some auxiliary lemmas needed in other parts of the paper.

**Acknowledgements.** We thank Ján Pich for answering questions about [PS21]. We are grateful to Anupam Das for a discussion on extracting computational content from proofs and to Junhua Yu for comments on the game-theoretic witnessing argument. We also thank Marco Carmosino, Emil Jeřábek, Valentine Kabanets, Antonina Kolokolova, and Jan Krajíček for related discussions. This work received support from the Royal Society University Research Fellowship URF\R1\191059 and from the EPSRC New Horizons Grant EP/V048201/1.

## 2 Preliminaries

This section presents some basic definitions and fixes notation.

### 2.1 Complexity theory

Given a function  $t: \mathbb{N} \rightarrow \mathbb{N}$ , we generalize the definition of each level of the polynomial hierarchy to machines that run in time  $t(n)$  in the natural way. For a fixed  $i \geq 1$ , we let  $\Pi_i\text{-TIME}[t]$  denote the set of languages  $L$  that admit a deterministic machine  $A$  running in time  $t(n)$  such that, for every  $x \in \{0, 1\}^n$ ,

$$x \in L \iff \forall z_1 \in \{0, 1\}^{\leq t(n)} \exists z_2 \in \{0, 1\}^{\leq t(n)} \dots Q_i z_i \in \{0, 1\}^{\leq t(n)} A(x, z_1, \dots, z_i) = 1.$$

The class  $\Sigma_i\text{-TIME}$  is defined in an analogous way. This generalises the classes  $\Sigma_i^p$  and  $\Pi_i^p$  corresponding to the  $i$ -th level of the polynomial hierarchy.

We consider (non-uniform) Boolean circuits over a standard set of gates of fan-in at most two, such as  $\{\wedge, \vee, \neg\}$ . The size of a circuit is the number of gates in the circuit. We adopt this convention only for concreteness, as our results are robust and do not depend on specific details of the circuit model. We let  $\text{SIZE}[s]$  denote the set of languages that admit non-uniform Boolean circuits of size  $s(n)$ .

We also consider circuits and corresponding circuit classes obtained by extending deterministic circuits to circuits with a constant number of alternations. For a fixed  $i \geq 1$ , we say that a language  $L \in \Sigma_i\text{-SIZE}[s]$  if there is a sequence  $\{C_n\}_{n \geq 1}$  of deterministic Boolean circuits  $C_n$  of size  $s(n)$  such that, for every  $x \in \{0, 1\}^n$ ,

$$x \in L \iff \exists z_1 \in \{0, 1\}^{s(n)} \forall z_2 \in \{0, 1\}^{s(n)} \dots Q_i z_i \in \{0, 1\}^{s(n)} C_n(x, z_1, \dots, z_i) = 1.$$

The class  $\Pi_i\text{-SIZE}[s]$  is defined in an analogous way. For convenience, we might refer to  $\Sigma_1\text{-SIZE}[s]$  as  $\text{NSIZE}[s]$ , i.e., the set of languages computed by non-deterministic circuits of size at most  $s(n)$ . When we write  $C(x) = 1$  for a non-deterministic circuit  $C$  and input  $x$ , we implicitly refer to its acceptance condition, i.e., that there is an input  $z$  such that  $C(x, z) = 1$ . We adopt the analogous convention for co-nondeterministic circuits and for circuit classes with additional alternations.

We will also consider languages that are computed by circuits with oracle gates. For an oracle  $\mathcal{O}$ , we let  $\text{SIZE}^{\mathcal{O}}[s]$  denote the set of languages computed by circuits of size at most  $s$  that can also make use of  $\mathcal{O}$ -oracle gates.

Finally, for convenience we often abuse notation and associate the size of a Boolean circuit to its bit-length, i.e., its description length under a reasonable encoding.

### 2.2 Logic and bounded arithmetic

We refer to [Bus97] for an introduction to bounded arithmetic and to the textbooks [Kra95, CN10] for a comprehensive treatment. Below we review the relevant definitions and fix notation.

We use  $\mathcal{L}(\mathcal{T})$  to denote the language (vocabulary) of a theory  $\mathcal{T}$ .

For a structure  $\mathcal{M}$  over a language  $\mathcal{L}$ , we often write  $\mathcal{M} = (\mathcal{D}, \mathcal{I})$  to explicitly refer to its domain  $\mathcal{D}$  and interpretations  $\mathcal{I}$ . As usual, the  $\mathcal{M}$ -interpretation of a function symbol  $f \in \mathcal{L}$  will be denoted by  $f^{\mathcal{M}}$  (similarly for relations and constants). The  $\mathcal{M}$ -interpretation of an  $\mathcal{L}$ -term  $t$  is also denoted by  $t^{\mathcal{M}}$ .

Given a formula  $\psi$ , we write  $\psi(y)$  to explicitly indicate that  $y$  may be a free variable in  $\psi$ . For a formula  $\varphi(x)$  and a term  $t$ , we write  $\varphi(x/t)$  for substitution of the free variable  $x$  with  $t$ , or simply  $\varphi(t)$  if it is clear from the context. Similarly, we use  $s(x)$  to denote a term  $s$  that may contain  $x$  as a free variable, and  $s(x/t)$

to denote the substitution of the free variable  $x$  with  $t$ , or simply  $s(t)$  if it is clear.

**The language  $\mathcal{L}_{PV}$ .** In theoretical computer science one typically considers functions and predicates that operate over binary strings. For the computational models considered in this paper, this is equivalent to operations on integers, by identifying each non-negative integer with its binary representation. For convenience, we adopt the latter perspective when introducing the language (vocabulary)  $\mathcal{L}_{PV}$  of theories  $T_{PV}^i$ .

Let  $\mathbb{N}$  denote the set of non-negative integers. For  $a \in \mathbb{N}$ , we let  $|a| = \max\{\lceil \log_2(a+1) \rceil, 1\}$  denote the length of the binary representation of  $a$ . For a constant  $k \geq 1$ , we say that a function  $f: \mathbb{N}^k \rightarrow \mathbb{N}$  is computable in polynomial time if  $f(x_1, \dots, x_k)$  can be computed in time polynomial in  $|x_1|, \dots, |x_k|$ . Recall that FP denotes the set of polynomial time functions. While this definition refers to a particular model of computation (Turing machines), Cobham [Cob65] proved that FP can be introduced in a machine independent way as the closure of a set of base functions under composition and limited recursion on notation. We briefly review this construction.<sup>5</sup>

Consider the following class  $\mathcal{F}_0$  of base functions:

$$c(x) = 0, \quad s_0(x) = 2 \cdot x, \quad s_1(x) = 2x + 1, \quad \pi_\ell^i(x_1, \dots, x_\ell) = x_i, \quad x \# y = 2^{|x| \cdot |y|}$$

We say that a function  $f(\vec{x}, y)$  is defined from functions  $g(\vec{x})$ ,  $h_0(\vec{x}, y, z)$ ,  $h_1(\vec{x}, y, z)$ , and  $k(\vec{x}, y)$  by *limited recursion on notation* if

$$\begin{aligned} f(\vec{x}, 0) &= g(\vec{x}) \\ f(\vec{x}, s_0(y)) &= h_0(\vec{x}, y, f(\vec{x}, y)) \\ f(\vec{x}, s_1(y)) &= h_1(\vec{x}, y, f(\vec{x}, y)) \\ f(\vec{x}, y) &\leq k(\vec{x}, y) \end{aligned}$$

for every sequence  $\vec{x}$  and  $y$  of natural numbers. Let  $\mathcal{F}$  be the least class of functions that contains  $\mathcal{F}_0$  and is closed under composition and limited recursion on notation. Cobham [Cob65] proved that  $f \in \mathcal{F}$  if and only if  $f \in \text{FP}$ .

We let  $\mathcal{L}_{PV}$  contain the constant symbols 0 and 1, and a function symbol  $f$  for every function in FP. In particular,  $\mathcal{L}_{PV}$  contains function symbols for the length function  $|x|$ ,  $\leq$ ,  $+$ , etc.<sup>6</sup>

We use the standard notation  $n \in \text{Log}$  and  $n \in \text{LogLog}$  for  $\exists N \ n = |N|$  and  $\exists N \ n = ||N||$ , respectively. In particular, we define  $\forall n \in \text{Log}$  and  $\forall n \in \text{LogLog}$  as  $\forall N \ \forall n = |N|$  and  $\forall N \ \forall n = ||N||$ , respectively.

**Bounded formulas and theories  $T_{PV}^i$ .** A *bounded quantifier* is a quantifier of the form  $Qx \leq t$ , where  $Q \in \{\exists, \forall\}$  and  $t$  is an  $\mathcal{L}_{PV}$ -term that does not involve  $x$ .<sup>7</sup> An  $\mathcal{L}_{PV}$ -formula  $\psi$  is *bounded* if every quantifier in  $\psi$  is bounded.

We will need to introduce a hierarchy of bounded formulas to define the theories  $T_{PV}^i$ . We let  $\Sigma_0^b = \Pi_0^b$  be the set of quantifier-free  $\mathcal{L}_{PV}$ -formulas. We then recursively define sets  $\Sigma_i^b$  and  $\Pi_i^b$  of formulas as follows.

<sup>5</sup>This is not strictly needed in our presentation. We include it here because it provides more intuition about the language of theories  $T_{PV}^i$  and the typical choice in bounded arithmetic of defining FP over non-negative integers instead of binary strings.

<sup>6</sup>It is also possible to include in  $\mathcal{L}_{PV}$  a function symbol for every polynomial time *algorithm*, where an algorithm can be described from the base functions and operations allowed in Cobham's characterisation. However, this is inessential. The theories  $T_{PV}^i$  will contain all true universal sentences, and polynomial time algorithms that compute the same function are provably equivalent in these theories.

<sup>7</sup>Bounded quantifiers can be expressed with the usual quantifiers from first-order logic. For instance, a formula  $\psi(y)$  of the form  $\forall x \leq t(y) \ \varphi(x, y)$  is equivalent to  $\forall x (x \leq t(y) \rightarrow \varphi(x, y))$ . On the other hand, a formula of the form  $\exists x \leq t(y) \ \varphi(x, y)$  is equivalent to  $\exists x (x \leq t(y) \wedge \varphi(x, y))$ .



For each  $i \geq 1$ ,  $\Sigma_i^b$  and  $\Pi_i^b$  constitute the smallest class of  $\mathcal{L}_{PV}$ -formulas such that the following conditions hold:

1.  $\Sigma_{i-1}^b \cup \Pi_{i-1}^b \subseteq \Sigma_i^b \cap \Pi_i^b$ ;
2. both  $\Sigma_i^b$  and  $\Pi_i^b$  are closed under Boolean connectives  $\wedge$  and  $\vee$ ;
3. if  $\psi(\vec{x}) \equiv \exists y \leq t(\vec{x}) \varphi(\vec{x}, y)$  is a bounded formula and  $\varphi \in \Sigma_i^b$ , then  $\psi \in \Sigma_i^b$ ;
4. similarly, if  $\psi(\vec{x}) \equiv \forall y \leq t(\vec{x}) \varphi(\vec{x}, y)$  is a bounded formula and  $\varphi \in \Pi_i^b$ , then  $\psi \in \Pi_i^b$ ;
5. the negation  $\neg\psi$  of a formula  $\psi$  from  $\Sigma_i^b$  is in  $\Pi_i^b$  and vice versa.

(For convenience, we sometimes describe formulas with the implication symbol  $\rightarrow$ , implicitly assuming that it is expressed using the Boolean connectives appearing above.)

Note that to each  $\mathcal{L}_{PV}$ -formula  $\phi(x_1, \dots, x_k)$  we can associate a language  $L_\phi \subseteq \{0, 1\}^*$  consisting of binary encodings of all tuples  $(a_1, \dots, a_k) \in \mathbb{N}^k$  such that  $\mathbb{N} \models \phi(a_1, \dots, a_k)$ . It is known that  $\phi \in \Sigma_i^b$  if and only if  $L_\phi \in \Sigma_i^p$  [Sto76, Wra76, KH82], where  $\Sigma_i^p$  denotes the  $i$ -th level of the polynomial hierarchy.

For  $j \geq 0$ , we let  $\forall\Sigma_j^b$  denote the set of  $\mathcal{L}_{PV}$ -sentences of the form  $\forall \vec{y} \varphi(\vec{y})$ , where  $\varphi$  is a  $\Sigma_j^b$ -formula. We sometimes write  $\Sigma_i^b(\mathcal{L})$ ,  $\Pi_i^b(\mathcal{L})$ , and  $\forall\Sigma_i^b(\mathcal{L})$  to emphasize the underlying language  $\mathcal{L}$  of a class of formulas.

As expected, the intended model of theories  $T_{PV}^i$  is  $\mathbb{N}$ , with the interpretation of each function symbol  $f \in \mathcal{L}_{PV}$  as the corresponding polynomial time function. We will refer to  $(\mathbb{N}, 0^{\mathbb{N}}, +^{\mathbb{N}}, \dots)$  as the *standard model*.

**Definition 2.1** (Theories  $T_{PV}^i$ ). For each integer  $i \geq 1$ , we let  $T_{PV}^i$  denote the theory of all true (with respect to  $\mathbb{N}$ )  $\forall\Sigma_{i-1}^b$  sentences over the language  $\mathcal{L}_{PV}$ .

Note that  $T_{PV}^1$  is the theory of true universal sentences. For simplicity, we might refer to  $T_{PV}^1$  just as  $T_{PV}$ .

In order to simplify the presentation of some results, we introduce the following definition.

**Definition 2.2** (Closure under if-then-else). A theory  $\mathcal{T}$  is *closed under if-then-else* if for every quantifier-free formula  $\varphi(\vec{x})$  and terms  $t_1(\vec{x})$  and  $t_2(\vec{x})$ , there exists a term  $t(\vec{x})$  such that

$$\mathcal{T} \vdash (t(\vec{x}) = t_1(\vec{x}) \wedge \varphi(\vec{x})) \vee (t(\vec{x}) = t_2(\vec{x}) \wedge \neg\varphi(\vec{x})).$$

We note that in such a theory the provability of a disjunction  $\psi(x, t_1(x)) \vee \psi(x, t_2(x)) \vee \dots \vee \psi(x, t_k(x))$  yields the provability of  $\psi(x, t(x))$ , for a quantifier-free formula  $\psi(x, y)$ . Typical theories of bounded arithmetic (e.g.,  $S_2^1$  and  $T_{PV}^i$ ) are closed under if-then-else or admit a suitable extension that is closed under this property.

**Theory APC<sub>1</sub>**. To formalize probabilistic methods and randomised algorithms, Jeřábek [Jeř07a] introduced the theory APC<sub>1</sub> by extending PV with the *dual Weak Pigeonhole Principle* of PV functions, i.e., there is no PV function  $f : [2^n] \rightarrow [(1 + 1/n) \cdot 2^n]$  that is surjective.<sup>8</sup> More formally, we define dWPHP( $f$ ) for a function  $f$  as the sentence<sup>9</sup>

$$\text{dWPHP}(f) \triangleq \forall n \in \text{Log} \ \forall \vec{z} \ \exists y < (1 + 1/n) \cdot 2^n \ \forall x < 2^n \ f(\vec{z}, x) \neq y. \quad (1)$$

<sup>8</sup>The size of the codomain (with respect to the size of the domain) affects the power of the dual Weak Pigeonhole Principle. This can be a subtle point, as the equivalence between dual Weak Pigeonhole Principles with different codomain sizes is not known to be provable in PV (see [Jeř07b] for more details).

<sup>9</sup>Note that the additional parameter  $\vec{z}$  is crucial in the definition of APC<sub>1</sub>. If we remove this parameter in the definition of dWPHP, denoted by dWPHP', the theory PV + dWPHP'(PV) will be a (possibly) weaker fragment of APC<sub>1</sub> (see, e.g., [PS21]).

Let  $\text{dWPHP}(\text{PV}) \triangleq \{\text{dWPHP}(f) \mid f \in \mathcal{L}_{\text{PV}}\}$ . Then  $\text{APC}_1 \triangleq \text{PV} + \text{dWPHP}(\text{PV})$ . (For a definition of theory PV, see [Kra95].)

Jeřábek [Jeř04, Jeř05, Jeř07a, Jer09] developed a complicated (but intuitive) framework for approximate counting in  $\text{APC}_1$  built on a formalisation of the Nisan-Wigderson PRG [NW94] in  $\text{APC}_1$ .

By counting the quantifier alternations in Equation (1), it is easy to see that  $\text{dWPHP}(f)$  is a  $\forall\Sigma_2^b$ -sentence in  $\mathcal{L}_{\text{PV}}$ . As a result,  $\text{APC}_1$  is a subtheory of  $\text{T}_{\text{PV}}^3$ . We note that our unprovability result for  $\text{APC}_1$  (Corollary 1.4) is quite robust and works with any non-trivial codomain size in the definition of  $\text{dWPHP}(f)$ , since this does not increase the quantifier complexity of the corresponding sentences.

### 3 Auxiliary Results in Logic and Complexity

This section describes auxiliary results that will be used to show the unprovability of strong circuit lower bounds in bounded arithmetic.

#### 3.1 Total search problems and the polynomial hierarchy

In this section, we define complexity classes and circuit classes associated with total search problems in the polynomial hierarchy and explore their basic properties.

Recall that a relation  $R \subseteq \{0, 1\}^* \times \{0, 1\}^*$  is in TFNP and if there is a polynomial  $p(n)$  and an efficient machine  $A$  such that

- For every  $x \in \{0, 1\}^*$  there is  $y \in \{0, 1\}^{\leq p(|x|)}$  such that  $R(x, y)$  holds. Moreover, any such  $y$  is of length at most  $p(|x|)$ .
- For every pair  $(x, y)$  of strings  $x, y \in \{0, 1\}^*$ ,  $(x, y) \in R$  if and only if  $A(x, y) = 1$ .

The next definition is a standard generalisation of this class.

**Definition 3.1.** For  $i \geq 1$ , we say that a relation  $R \in \text{TF}\Sigma_i^p$  if there is a polynomial  $p(n)$  and an efficient machine  $A$  such that the following conditions hold:

- For every  $x \in \{0, 1\}^*$  there is  $y \in \{0, 1\}^{\leq p(|x|)}$  such that  $R(x, y)$  holds.
- For every pair  $(x, y)$  of strings  $x, y \in \{0, 1\}^*$ ,

$$R(x, y) \iff \forall z_1 \in \{0, 1\}^{p(|x|)} \exists z_2 \in \{0, 1\}^{p(|x|)} \dots Q_{i-1} z_{i-1} \in \{0, 1\}^{p(|x|)} A(x, y, z_1, \dots, z_{i-1}).$$

In other words,  $R \in \Pi_{i-1}^p$ .

We will need the following simulation result.

**Theorem 3.2.** For every  $i \geq 1$  and  $s(n) \geq n$ ,  $\text{SIZE}^{\Sigma_{i-1}^p}[s(n)] \subseteq \Sigma_i\text{-SIZE}[\text{poly}(s(n))]$ .

*Proof.* The proof is similar to the well-known inclusion  $\text{P}^{\Sigma_{i-1}^p} \subseteq \Sigma_i^p$  (see, e.g., [Pap94, Chapter 17]), and we omit the details.  $\square$

#### 3.2 The Nisan-Wigderson generator

In this section, we review basic properties of the Nisan-Wigderson [NW94] pseudorandom generator and fix notation. For an introduction to this generator and to computational pseudorandomness, see [Vad12].

**Definition 3.3.** A collection  $\mathcal{S} = \{S_1, \dots, S_k\}$  of sets  $S_i$  is said to be an  $(m, \ell, a)$ -design if

- for every  $i \in [k]$ ,  $S_i \subseteq [m]$ ;
- for every  $i \in [k]$ ,  $|S_i| = \ell$ ; and
- for every  $i \neq j \in [k]$ ,  $|S_i \cap S_j| \leq a$ .<sup>10</sup>

The *size* of a design  $\mathcal{S}$  is defined as the number of sets in  $\mathcal{S}$ .

**Lemma 3.4** (Explicit designs; see, e.g., [NW94, Vad12]). *For every constant  $c \geq 2$  and every sufficiently large  $n \in \mathbb{N}$ , there exists an  $(n^c, n^{c/2}, n)$ -design  $\mathcal{S}_{c,n}$  of size  $2^n$ . Moreover, for every fixed  $c$ , there is an algorithm that, given a large enough  $n$  and an index  $i \in [2^n]$ , outputs the  $i$ -th set  $S_i \in \mathcal{S}_{c,n}$  in time  $\text{poly}(n)$ .*

Recall that, given an  $(m, \ell, a)$ -design  $\mathcal{S}$  of size  $N$  and a function  $f: \{0, 1\}^\ell \rightarrow \{0, 1\}$ , the *Nisan-Wigderson generator* (NW generator) maps a *seed*  $w \in \{0, 1\}^m$  into the  $N$ -bit string

$$f(w|_{S_1})f(w|_{S_2}) \dots f(w|_{S_N}),$$

where  $w|_{S_i}$  is the string of length  $\ell$  obtained from  $w$  by selecting the bits indexed by  $S_i \in \mathcal{S}$ .

It will be convenient to view the NW generator as a Boolean function and to introduce additional notation. For a large constant  $c \geq 1$ , given a function  $h: \{0, 1\}^{n^{c/2}} \rightarrow \{0, 1\}$ , we will use the NW generator to define a function  $\text{NW}_h: \{0, 1\}^{n^c} \times \{0, 1\}^n \rightarrow \{0, 1\}$ . More precisely,

- The seed length is  $n^c$ .
- The corresponding design is described by a  $2^n \times n^c$  Boolean matrix  $A$  where each row has exactly  $n^{c/2}$  entries set to 1, and the 1 entries in distinct rows overlap in at most  $n$  columns. As stated in Lemma 3.4, designs with these parameters are known to exist. Given a pair  $(i, j) \in [2^n] \times [n^c]$ , the  $(i, j)$ -entry of the corresponding design matrix can be explicitly computed by circuits of size  $\text{poly}(n)$  [NW94].
- For a row index  $x \in \{0, 1\}^n$  of  $A$ , we use  $J_x \subseteq [n^c]$  to denote the set of  $n^{c/2}$  columns of the  $x$ -th row of  $A$  set to 1.
- It will often be convenient to consider an  $n^c$ -bit string  $w$  as a function in  $\{0, 1\}^{[n^c]}$  that maps  $[n^c]$  to  $\{0, 1\}$ . If  $S_1, S_2 \subseteq [n^c]$  is a partition of  $[n^c]$ ,  $a \in \{0, 1\}^{S_1}$ , and  $u \in \{0, 1\}^{S_2}$ , we let  $w = u \cup a$  denote the corresponding  $n^c$ -bit string obtained from  $a$  and  $u$ .<sup>11</sup>
- For  $x \in \{0, 1\}^n$  and strings  $a \in \{0, 1\}^{n^c - n^{c/2}}$  and  $u \in \{0, 1\}^{n^{c/2}}$ , we let  $r_x(a, u)$  denote the string  $w = u \cup a$  of length  $n^c$  obtained by viewing  $a \in \{0, 1\}^{[n^c] \setminus J_x}$  and  $u \in \{0, 1\}^{J_x}$ .
- By fixing the seed  $w \in \{0, 1\}^{n^c}$  in the NW generator and the function  $h: \{0, 1\}^{n^{c/2}} \rightarrow \{0, 1\}$ , we obtain a function  $\text{NW}_h(w): \{0, 1\}^n \rightarrow \{0, 1\}$  in the natural way. Similarly, we can obtain a family  $\{\text{NW}_h(w)\}_{w \in \{0, 1\}^{n^c}}$  of functions, one for each possible seed  $w$ .

### 3.3 Hardness amplification in the polynomial hierarchy

In order to relax the average-case complexity parameter in our unprovability results, we prove a hardness amplification theorem for the polynomial hierarchy. This result can be seen as the “relativised” version of [HVV06] (see also [PS21, Section 3.3]). For completeness, we sketch their proof and explain how to adapt their result for our purpose.

<sup>10</sup>Designs are also called combinatorial designs by some authors. We will use both terms interchangeably.

<sup>11</sup>This notation is consistent with the standard set-theoretic definition of a function as a set of pairs.

**Theorem 3.5.** *There is a constant  $\gamma > 0$  and  $\ell = \ell(n) = \text{poly}(n)$  such that the following holds for every  $i \geq 1$ . Let  $s_1, s_2: \mathbb{N} \rightarrow \mathbb{N}$  be non-decreasing functions, where  $s_2(n) = n^{\omega(1)}$ , and suppose there is a function  $f_n: \{0, 1\}^n \rightarrow \{0, 1\}$  computable by  $\Sigma_i\text{-SIZE}[s_1(n)]$  circuits (resp.  $\Pi_i\text{-SIZE}[s_1(n)]$  circuits) such that each  $\Sigma_{i-1}^p$ -oracle circuit  $A_n$  of size at most  $s_2(n)$  satisfies*

$$\Pr_{x \in \{0,1\}^n} [f_n(x) = A_n(x)] \leq 1 - \frac{1}{n}.$$

*Then there exist a function  $h_\ell: \{0, 1\}^\ell \rightarrow \{0, 1\}$  computable by  $\Sigma_i\text{-SIZE}[\text{poly}(\ell) \cdot s_1(\ell)]$  circuits (resp.  $\Pi_i\text{-SIZE}[\text{poly}(\ell) \cdot s_1(\ell^\gamma)]$  circuits) such that each  $\Sigma_{i-1}^p$ -oracle circuit  $B_\ell$  of size at most  $s_2(\ell^\gamma)^\gamma$  satisfies*

$$\Pr_{y \in \{0,1\}^\ell} [h_\ell(y) = B_\ell(y)] \leq \frac{1}{2} + \frac{1}{s_2(\ell^\gamma)^\gamma}.$$

Note that since  $\Sigma_{i-1}^p$ -oracle circuits are closed under complementation, we only need to prove the case where  $f_n$  is computable by  $\Sigma_i$  circuits. More formally, given a function  $f: \{0, 1\}^n \rightarrow \{0, 1\}$  computable by  $\Pi_i\text{-SIZE}[s_1(n)]$  circuits that is hard on average against  $\Sigma_{i-1}^p$ -oracle circuits of size  $s_2(n)$ , we can consider  $g(x) \triangleq \neg f(x)$  that is computable by  $\Sigma_i\text{-SIZE}[s_1(n)]$  circuits and still hard on average against the same class. By the hardness amplification theorem for  $\Sigma_i$  circuits, we can obtain a function  $h_\ell$  computable by  $\Sigma_i\text{-SIZE}[\text{poly}(\ell) \cdot s_1(\ell)]$  circuits that is strongly hard on average against  $\Sigma_{i-1}^p$ -oracle circuits of size  $s_2(\ell^\gamma)^\gamma$ . The negation of  $h_\ell$  is then the required hard function computable in  $\Pi_i\text{-SIZE}[\text{poly}(\ell) \cdot s_1(\ell)]$ .

In the rest of the section we will sketch the proof of Theorem 3.5. We first fix the notation.

- A *probabilistic function* is a Boolean function with two inputs  $h(x; r)$  where the second input is treated as random bits. If the random bits are omitted, a probabilistic function is treated as a function mapping the input to a random variable distributed according to the output of the function over the random bits.
- Let  $g$  be a function probabilistic function) with input length  $n$ , the  $k$ -th *direct product* is defined as the function (resp. probabilistic function) with input length  $k \cdot n$  and output length  $k$  as follows:

$$g^{\otimes k}(x_1, \dots, x_k) \triangleq g(x_1) \parallel \dots \parallel g(x_k).$$

- The *bias* of a random variable  $X$  is defined as  $\text{Bias}(X) \triangleq |\Pr[X = 0] - \Pr[X = 1]|$ . The *bias* of a probabilistic function  $h(x; r)$  is defined as the bias of the random variable  $h(x; r)$  for a uniformly random  $x$  and  $r$ . The probabilistic function  $h$  is said to be *balanced* if  $\text{Bias}(h) = 0$ .
- A probabilistic function  $h: \{0, 1\}^n \times \{0, 1\}^r \rightarrow \{0, 1\}$  is  $\delta$ -*random* if  $h$  is balanced and there is a subset  $H \subseteq \{0, 1\}^n$  of size  $2\delta \cdot 2^n$  such that  $h$  is a “coin flip” over  $H$  and deterministic outside  $H$  (i.e.,  $\Pr[h(x) = 1] = 1/2$  for every  $x \in H$ , and  $h(x)$  is deterministic for every  $x \notin H$ ).
- The *expected bias* of a probabilistic function  $h$  is defined as  $\text{ExpBias}(h) \triangleq \mathbb{E}_x [\text{Bias}(h(x))]$ .
- The *noise stability* of a Boolean function  $C: \{0, 1\}^k \rightarrow \{0, 1\}$  with respect to the noise rate  $\delta$  is defined as

$$\text{NoiseStab}_\delta(g) \triangleq 2 \cdot \Pr_{x, \eta} [C(x) = C(x \oplus \eta)] - 1,$$

where  $x \sim \{0, 1\}^k$  and each bit of  $\eta$  is 1 independently with probability  $\delta$ . By Lemma 3.7 of [HVV06],  $\text{ExpBias}[C \circ g^{\otimes k}] \leq \sqrt{\text{NoiseStab}_\delta[C]}$  for every  $\delta$ -random probabilistic function  $g$ .

- Two random variables  $X_1$  and  $X_2$  are said to be  $\varepsilon$ -indistinguishable for size  $s$ , denoted by  $X_1 \approx_\varepsilon^s X_2$ , if for every  $\Sigma_{i-1}^p$ -oracle circuit  $C$  of size  $s$ ,  $|\Pr[C(X_1) = 1] - \Pr[C(X_2) = 1]| \leq \varepsilon$ . Note that our definition of the indistinguishability differs from the original definition in [HVV06] since we are proving hardness amplification against  $\Sigma_{i-1}^p$ -oracle circuits.
- For simplicity, we say a function  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  is  $\varepsilon$ -hard for size  $s$ , if for every  $\Sigma_{i-1}^p$ -oracle circuit  $C$  of size  $s$ ,  $C(x) = f(x)$  for at most an  $\varepsilon$  fraction of  $x \in \{0, 1\}^n$ .

We assume that  $f_n$  is *balanced*, that is,  $\Pr_x[f_n(x) = 1] = 1/2$  for every  $n \geq 1$ . This is without loss of generality, since we can first increase the input length by one then use non-uniformity to make the resulting function balanced, without a relevant change of parameters.

The hardness amplification of [HVV06] proceeds as follows. (We discuss the proofs of the lemmas stated below in Appendix D.)

**The Construction.** Fix any  $n \geq 1$ . Let  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  be the hard function and  $C : \{0, 1\}^k \rightarrow \{0, 1\}$  be an explicit circuit to be determined later. Let  $G : \{0, 1\}^\ell \rightarrow (\{0, 1\}^n)^k$  be an explicit function in the sense that given  $\sigma \in \{0, 1\}^\ell$  and  $i \in [k]$ , we can compute the  $i$ -th block  $X_i \in \{0, 1\}^n$  of the output of  $G(\sigma)$  in  $\text{poly}(\ell, \log k)$  time.<sup>12</sup> The amplified function is defined as  $\text{Amp}_f : \{0, 1\}^\ell \rightarrow \{0, 1\}$ :

$$\text{Amp}_f(\sigma) \triangleq C(f(X_1), f(X_2), \dots, f(X_k)),$$

where  $(X_1, X_2, \dots, X_k) \triangleq G(\sigma)$ . We need to carefully choose  $C$  and  $G$  such that  $\text{Amp}_f(\sigma)$  is computable in  $\Sigma_i\text{-SIZE}[\text{poly}(n) \cdot s_1(n)]$  and can amplify the hardness of  $f$ .

**The Choice of  $G$ .** To ensure the hardness of  $\text{Amp}_f$ , the function  $G_k : \{0, 1\}^\ell \rightarrow (\{0, 1\}^n)^k$  should satisfy the following two technical requirements.

- $G_k$  is *indistinguishability-preserving for size  $t = k^2$* : Let  $f_1, \dots, f_k, g_1, \dots, g_k$  be probabilistic functions such that for every  $i \in [k]$ ,  $x \| f_i(x) \approx_\varepsilon^s x \| g_i(x)$  for  $x \sim \{0, 1\}^n$ , then

$$\sigma \| f_1(X_1) \| \dots \| f_k(X_k) \approx_{k \cdot \varepsilon}^{s-t} \sigma \| g_1(X_1) \| \dots \| g_k(X_k),$$

where  $\sigma \sim \{0, 1\}^\ell$  and  $(X_1, \dots, X_k) \triangleq G_k(\sigma)$ .

- $G_k$  is  $2^{-n}$ -pseudorandom against (read-once oblivious) branching programs of size  $2^n$  and block-size  $n$ :<sup>13</sup> for every branching program  $B$  of size  $2^n$  and block-size  $n$ , we have

$$\left| \Pr_{x \sim \{0, 1\}^\ell} [B(G_k(x)) = 1] - \Pr_{y \sim \{0, 1\}^{nk}} [B(y) = 1] \right| \leq 2^{-n}.$$

**Lemma 3.6** (Generalized version of [HVV06, Theorem 5.12]). *For every  $k \leq 2^n$ , there is an explicitly computable generator  $G_k$  (in the sense that the  $i$ -th block of  $G_k(\sigma)$  can be computed in  $\text{poly}(\ell, \log k)$  time) that satisfies the requirements above with seed length  $\ell = O(n^2)$ .*

<sup>12</sup>We note that  $C$  is used to replace the XOR function in the standard hardness amplification based on Yao's XOR Lemma (see, e.g., Theorem 19.2 of [AB09]), while  $G$  is used as a pseudorandom generator that (in some sense) "fools"  $C \circ f^{\otimes k}$ .

<sup>13</sup>See [HVV06, Definition 5.4] for the precise definition of this branching program model.

**The Choice of  $C$ .** The outer function  $C$ , which serves as the counterpart of the XOR function in Yao's XOR Lemma (see, e.g., [AB09, Theorem 12.9]), is chosen according to the following lemma.

**Lemma 3.7** (Generalized version of [HVV06, Lemma 5.15]). *For every  $i \geq 1$ ,  $\delta(n) = 1/\text{poly}(n)$ , and  $k = k(n)$  such that  $n^{\omega(1)} \leq k \leq 2^n$ , there is a function  $C_k : \{0, 1\}^k \rightarrow \{0, 1\}$  such that:*

- (i)  $\text{NoiseStab}_\delta[C_k] \leq 1/k^{\Omega(1)}$ ;
- (ii) *For every  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  computable by  $\Sigma_i\text{-SIZE}[s(n)]$  circuits,  $(C_k \circ f^{\otimes k}) \circ G_k : \{0, 1\}^\ell \rightarrow \{0, 1\}$  is computable by  $\Sigma_i\text{-SIZE}[\text{poly}(n) \cdot s(n)]$  circuits.*
- (iii)  $C_k$  is computable by a branching program of size  $\text{poly}(n) \cdot k$  and by a deterministic circuit of size  $\text{poly}(n) \cdot k$ .

Note that the second item means that the function  $(C_k \circ f^{\otimes k}) \circ G$  is efficiently computable *even if  $k$  is as large as  $2^n$* . The argument relies on the explicitness of  $C_k$  and  $G$  as well as on the power of  $\Sigma_i$ -circuits. This is crucial for hardness amplification up to  $1/2 - 1/s_2(\ell^\gamma)^\gamma$  (instead of only  $1/2 - 1/\text{poly}(\ell)$ ).

**Proof of the Hardness Amplification.** Following [HVV06, Section 5], we now argue that if  $f$  is  $\delta$ -hard for size  $s(n) \geq n^{\omega(1)}$ , where  $\delta \geq 1/\text{poly}(n)$ , then we can construct  $\text{Amp}_f : \{0, 1\}^\ell \rightarrow \{0, 1\}$  with  $\ell = \text{poly}(n)$  that is  $(1/2 - 1/s(\sqrt{\ell})^{\Omega(1)})$ -hard for size  $s(\sqrt{\ell})^{\Omega(1)}$ . To prove this, we need the following two technical lemmas.

**Lemma 3.8** ([HVV06, Lemma 5.7 and Lemma 5.12]). *Let  $g$  be an  $n$ -input single output  $\delta$ -random function, and  $C_k$  and  $G_k$  be defined as above. Then*

$$\text{ExpBias}[(C_k \circ g^{\otimes k}) \circ G_k] \leq \sqrt{\text{NoiseStab}_\delta(C_k) + 2^{-n+1}}.$$

**Lemma 3.9** (Generalized version of [HVV06, Lemma 5.2]). *Assume that  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  is  $\delta$ -hard for size  $s = n^{\omega(1)}$ . There is a  $\delta'$ -random function  $g$  with  $\delta' \in [\delta/2, \delta]$  such that  $\text{Amp}_f : \{0, 1\}^\ell \rightarrow \{0, 1\}$  has hardness*

$$\frac{1}{2} - \frac{\text{ExpBias}[(C \circ g^{\otimes k}) \circ G]}{2} - \frac{k}{s^{1/3}}$$

for size  $\Omega(s^{1/3}/\log(s/\delta)) - k^2 - \text{poly}(n) \cdot k$ .

Let  $k = k(n) = s(n)^{1/7}$ ,  $C_k$  be the function in Lemma 3.7, and  $G_k$  be the generator in Lemma 3.6 with  $\ell = O(n^2)$ . Recall that  $\text{Amp}_f : \{0, 1\}^\ell \rightarrow \{0, 1\}$  is defined as  $\text{Amp}_f \triangleq (C_k \circ f^{\otimes k}) \circ G_k$ , and note that the upper bound on the complexity of  $\text{Amp}_f$  is guaranteed by Lemma 3.7. By Lemma 3.9, we know that  $\text{Amp}_f$  has hardness

$$\frac{1}{2} - \frac{\text{ExpBias}[(C_k \circ g^{\otimes k}) \circ G_k]}{2} - \frac{k}{s^{1/3}} \quad (2)$$

for size  $\Omega(s^{1/3}/\log(s/\delta)) - k^2 - \text{poly}(n) \cdot k = s(\sqrt{\ell})^{\Omega(1)}$ , where  $g$  is some  $\delta'$ -random function with  $\delta' \in [\delta/2, \delta]$ . By Lemma 3.8, we can bound Equation (2) using the noise stability bound for  $C_k$  given in Lemma 3.7:

$$\begin{aligned} (2) &\geq \frac{1}{2} - \frac{\sqrt{\text{NoiseStab}_{\delta'}(C_k) + 2^{-n+1}}}{2} - \frac{k}{s^{1/3}} \\ &\geq \frac{1}{2} - \frac{\sqrt{k^{-\Omega(1)} + 2^{-n+1}}}{2} - \frac{k}{s(n)^{1/3}} \\ &\geq \frac{1}{2} - \frac{1}{s(\sqrt{\ell})^{\Omega(1)}}. \end{aligned}$$

This completes the argument. We discuss the proofs of the relevant technical lemmas in Appendix D.



### 3.4 Herbrand's Theorem and the KPT Witnessing Theorem

In this section, we review standard witnessing theorems previously used to show unprovability results in bounded arithmetic (see, e.g., [CKKO21, PS21]). In all results, we consider a universal theory  $\mathcal{T}$  with vocabulary  $\mathcal{L}$ .<sup>14</sup> (As a concrete example, one can take  $\mathcal{T} = \text{PV}$  and  $\mathcal{L} = \mathcal{L}_{\text{PV}}$ .)

**Two quantifiers ( $\forall\exists$ ).** The well-known Herbrand's theorem is the simplest witnessing result and can be applied to  $\forall\exists$ -sentences (see, e.g., Section 2 of [Koh08]).

**Theorem 3.10** (Herbrand's Theorem). *Let  $\mathcal{T}$  be a universal theory with vocabulary  $\mathcal{L}$ . If  $\mathcal{T} \vdash \forall x \exists y \varphi(x, y)$  for a quantifier-free  $\mathcal{L}$ -formula  $\varphi$ , there exist a constant  $\ell \geq 1$  and a sequence  $t_1, t_2, \dots, t_\ell$  of  $\mathcal{L}$ -terms such that*

$$\mathcal{T} \vdash \forall x (\varphi(x, t_1(x)) \vee \varphi(x, t_2(x)) \vee \dots \vee \varphi(x, t_\ell(x))).$$

*In particular, if  $\mathcal{T}$  is closed under if-then-else, then there is an  $\mathcal{L}$ -term  $t$  such that  $\mathcal{T} \vdash \forall x \varphi(x, t(x))$ .*

**Three quantifiers ( $\forall\exists\forall$ ).** The KPT Witnessing Theorem [KPT91] extends Herbrand's Theorem by providing witnessing functions for the existential quantifier in a provable  $\forall\exists\forall$ -sentence.

**Theorem 3.11** (KPT Witnessing [KPT91]). *Let  $\mathcal{T}$  be a universal theory with vocabulary  $\mathcal{L}$ . Suppose that, for a quantifier-free  $\mathcal{L}$ -formula  $\varphi$ ,  $\mathcal{T} \vdash \forall x \exists y \forall z \varphi(x, y, z)$ . Then there exist a constant  $\ell \geq 1$  and a sequence  $t_1, \dots, t_\ell$  of  $\mathcal{L}$ -terms such that*

$$\mathcal{T} \vdash \forall x \forall \vec{z} (\varphi(x, t_1(x), z_1) \vee \varphi(x, t_2(x, z_1), z_2) \vee \dots \vee \varphi(x, t_\ell(x, z_1, \dots, z_{\ell-1}), z_\ell)).$$

KPT witnessing has a well-known computational interpretation as an interactive game between a student and a teacher (see, e.g., [Pic15a]). In the first round, the student is given an arbitrary input  $x$ , and computes according to the term  $t_1(x)$ . This computation provides a candidate object  $y_1$ . The teacher then replies with an arbitrary “counterexample”  $z_1$  such that  $\neg\varphi(x, y_1, z_1)$  holds, whenever such  $z_1$  exists. Note that the next move of the student takes into account previously presented counterexamples, i.e., the term  $t_2$  depends on both  $x$  and  $z_1$ . According to Theorem 3.11, the game ends in at most  $\ell$  rounds, and the student is guaranteed to succeed, i.e., to output  $y$  such that  $\varphi(x, y, z)$  holds for every  $z$ .

*Example 3.12.* An example of the interactive game is the proof of the existence of two irrational numbers  $x, y$  such that  $x^y$  is rational (see Example 1.9), formalized (in some appropriate theory for real numbers) as:

$$\begin{aligned} &\exists x, y \in \mathbb{R} \exists p, q \in \mathbb{Z} \forall p', q' \in \mathbb{Z} \psi(x, y, p, q, p', q'), \text{ where} \\ &\psi(x, y, p, q, p', q') \triangleq x^y = p/q \wedge x \neq p'/q' \wedge y \neq p'/q' \end{aligned}$$

The student wants to learn  $x, y, p, q$  such that  $\psi(x, y, p, q, p', q')$  holds for every  $p', q'$ , with the help of a teacher that finds counterexamples  $p', q'$  making  $\psi(x, y, p, q, p', q')$  false. The student's strategy (say, extracted from the proof using KPT witnessing) is that:

- In the first round, try  $x = (\sqrt{2})^{\sqrt{2}}, y = \sqrt{2}, p = 2, q = 1$ , and ask for a counterexample  $p', q'$  from the teacher if it failed.
- Since  $y \neq p'/q'$ , the student knows that  $x = p'/q'$ . The student can then propose in the second round that  $x = \sqrt{2}, y = \sqrt{2}, p = p', q = q'$ .

<sup>14</sup>Recall that a theory  $\mathcal{T}$  is *universal* if all its axioms are universal formulas, i.e., a formula of the form  $\forall \vec{x} \varphi(\vec{x})$ , where  $\varphi$  is free of quantifiers.

**Four quantifiers ( $\forall\exists\forall\exists$ ).** It is also known that one can prove a witnessing theorem for  $\forall\exists\forall\exists$ -sentences using the standard model-theoretical proof of the KPT witnessing theorem.

**Theorem 3.13** (KPT Witnessing for  $\forall\exists\forall\exists$ -Sentences [KPT91]). *Let  $\mathcal{T}$  be a universal theory with vocabulary  $\mathcal{L}$ . Let  $\varphi$  be a quantifier-free  $\mathcal{L}$ -formula, and suppose that  $\mathcal{T} \vdash \forall x \exists y \forall z \exists w \varphi(x, y, z, w)$ . Then there is an  $\ell \geq 1$  and a finite sequence  $t_1, \dots, t_\ell$  of  $\mathcal{L}$ -terms such that*

$$\mathcal{T} \vdash \forall x, z_1, \dots, z_k \left( \psi(z, t_1(z), z_1) \vee \psi(x, t_2(x, z_1), z_2) \vee \dots \vee \psi(x, t_\ell(z_1, \dots, z_{\ell-1}), z_\ell) \right),$$

where  $\psi(x, y, z) \triangleq \exists w \varphi(x, y, z, w)$ .

For completeness, we provide a self-contained presentation of this result in Appendix B.

**Five or more quantifiers?** Unlike the case of four quantifiers, there is no obvious direct generalization of the KPT witnessing theorem to five or more quantifiers. The intuitive reason is that there is more than one universal quantifier within the outermost existential quantifier that we would like to witness, so the interaction pattern of the student and the teacher, which can provide counterexamples for all but the outermost universal quantifier, becomes much more complicated. Technically speaking, a close inspection of the model-theoretical proof of Theorem 3.13 presented in Appendix B reveals that the same argument cannot be applied to higher quantifier complexity.

### 3.5 A universal theory for $\mathcal{T}_{\text{PV}}^i$

There are two immediate issues when trying to show the unprovability of the lower bound sentence  $\text{LB}^i$  in  $\mathcal{T}_{\text{PV}}^i$ . Firstly,  $\text{LB}^i$  contains more quantifier alternations than a typical witnessing theorem can handle (see Section 4). Secondly,  $\mathcal{T}_{\text{PV}}^i$  is not a universal theory if  $i > 1$ , which violates a common assumption in these results. To address the latter, the first step of our argument is to turn  $\mathcal{T}_{\text{PV}}^i$  into a universal theory by introducing Skolem functions. In turn, this will allow us to reduce the quantifier complexity of  $\text{LB}^i$  so that the techniques developed in Section 4 can be applied (see Section 6.2).

**Theory  $\mathcal{U}_{\text{PV}}^i$  and Language  $\mathcal{L}_{\text{PV}}^i$ .** Let  $i \geq 1$ . For each  $(\Pi_{i-1}^b \cup \Sigma_{i-1}^b)$ -formula  $\alpha(\vec{z})$  over  $\mathcal{L}_{\text{PV}}$ , we introduce a function symbol  $f_\alpha$  interpreted (in the standard model) as the Boolean-valued function

$$f_\alpha^{\mathbb{N}}(\vec{z}) = \begin{cases} 1 & \text{if } \alpha^{\mathbb{N}}(\vec{z}) \text{ holds;} \\ 0 & \text{otherwise.} \end{cases}$$

Furthermore, when  $i \geq 2$ , for each  $\Sigma_{i-1}^b$ -formula  $\beta(\vec{x}, y)$  and term  $t$  in  $\mathcal{L}_{\text{PV}}$ , we introduce a function symbol  $g_{\beta, t}$  that is interpreted (in the standard model) as the function<sup>15</sup>

$$g_{\beta, t}^{\mathbb{N}}(\vec{x}) = \begin{cases} \text{smallest } y \in \mathbb{N} \text{ s.t. } \beta^{\mathbb{N}}(\vec{x}, y) & \text{if } \exists y \leq t^{\mathbb{N}}(\vec{x}) \beta^{\mathbb{N}}(\vec{x}, y); \\ 0 & \text{otherwise.} \end{cases}$$

Denote by  $\mathcal{L}_{\text{PV}}^i$  the language of  $\mathcal{L}_{\text{PV}}$  supplemented with the new function symbols. Let  $\mathcal{U}_{\text{PV}}^i$  be the theory consisting of all universal true sentences (over the standard model) in  $\mathcal{L}_{\text{PV}}^i$ .

<sup>15</sup>If the reader is somewhat uncomfortable with the possibility that the smallest  $y$  satisfying the condition below might be 0, we stress that this is not going to be an issue in our construction – see, e.g., the statement of Lemma 3.16.

**Correctness of the extension  $U_{PV}^i$ .** Now we show that  $U_{PV}^i$  is an extension of  $T_{PV}^i$ , that is, every sentence provable in  $T_{PV}^i$  is also provable in  $U_{PV}^i$ . First, we establish a few useful lemmas.

**Lemma 3.14.** *Let  $i \geq 0$ . For every  $\Sigma_i^b$ -formula (resp.  $\Pi_i^b$ -formula)  $\alpha(\vec{z})$  in the language  $\mathcal{L}_{PV}$ , there exists a formula  $\alpha^{\text{norm}}(\vec{z}) = Q_1x_1 \leq t_1(\vec{z}) \ Q_2x_2 \leq t_2(\vec{z}) \ \dots \ Q_ix_i \leq t_i(\vec{z}) \ \phi(\vec{z}, \vec{x})$ , where  $Q_1 = \exists$  (resp.  $Q_1 = \forall$ ),  $Q_j \in \{\forall, \exists\}$ , and  $Q_j \neq Q_{j+1}$  for every  $j \leq i-1$ , such that  $T_{PV}^1 \vdash \forall \vec{z} (\alpha(\vec{z}) \leftrightarrow \alpha^{\text{norm}}(\vec{z}))$ .*

*Proof.* Note that for every  $\Sigma_i^b$ -formula (resp.  $\Pi_i^b$ -formula)  $\alpha(\vec{z})$ , we can firstly find its prenex normal form  $\alpha^{\text{pnf}}(\vec{z})$  with  $i-1$  quantifier alternations starting with an existential (resp. universal) quantifier that is logically equivalent to  $\alpha(\vec{z})$ . Note that  $T_{PV}^1$  defines pairing and unpairing functions. Concretely, there are functions  $\langle \cdot, \cdot \rangle, \pi_1(\cdot), \pi_2(\cdot)$  such that  $T_{PV}^1 \vdash \forall x \forall y (\pi_1(\langle x, y \rangle) = x \wedge \pi_2(\langle x, y \rangle) = y \wedge |\langle x, y \rangle| \leq 10 \cdot (|x| + |y|))$ . It is then easy to see that

$$\begin{aligned} T_{PV}^1 \vdash & \left( \forall x \leq s \ \forall y \leq t \ \varphi(x, y) \right) \leftrightarrow \left( \forall p \leq (s \cdot t)^{10} (\pi_1(p) \leq s \wedge \pi_2(p) \leq t \rightarrow \varphi(\pi_1(p), \pi_2(p))) \right) \\ T_{PV}^1 \vdash & \left( \exists x \leq s \ \exists y \leq t \ \varphi(x, y) \right) \leftrightarrow \left( \exists p \leq (s \cdot t)^{10} (\pi_1(p) \leq s \wedge \pi_2(p) \leq t \wedge \varphi(\pi_1(p), \pi_2(p))) \right). \end{aligned}$$

Therefore we can further collapse adjacent quantifiers of the same kind to obtain  $\alpha^{\text{norm}}(\vec{z})$  as described above such that  $T_{PV}^1 \vdash \forall \vec{z} (\alpha(\vec{z}) \leftrightarrow \alpha^{\text{pnf}}(\vec{z}) \leftrightarrow \alpha^{\text{norm}}(\vec{z}))$ .  $\square$

**Lemma 3.15.** *For every  $i \geq 1$  and  $(\Pi_{i-1}^b \cup \Sigma_{i-1}^b)$ -formula  $\alpha(\vec{x})$ , we have (1)  $U_{PV}^i \vdash \forall \vec{x} (f_\alpha(\vec{x}) = 1 \leftrightarrow f_{\neg\alpha}(\vec{x}) = 0)$  and (2)  $U_{PV}^i \vdash \forall \vec{x} (f_\alpha(\vec{x}) = 0 \vee f_\alpha(\vec{x}) = 1)$ .*

*Proof.* By the definition of each  $f_\alpha^{\mathbb{N}}$ , we can see that the universal sentences  $\forall \vec{x} (f_\alpha(\vec{x}) = 1 \leftrightarrow f_{\neg\alpha}(\vec{x}) = 0)$  and  $\forall \vec{x} (f_\alpha(\vec{x}) = 0 \vee f_\alpha(\vec{x}) = 1)$  are both true in the standard model, so they are provable in  $U_{PV}^i$ .  $\square$

**Lemma 3.16** (Defining Axioms of  $g_{\beta,t}$ ). *Let  $i \geq 2$ ,  $\beta(\vec{x}, y)$  be any  $\Sigma_{i-1}^b$ -formula in  $\mathcal{L}_{PV}$ , and  $t$  be any term in  $\mathcal{L}_{PV}$ . Then  $U_{PV}^i \vdash \forall \vec{x} ((\exists y \leq t(\vec{x}) \ f_\beta(\vec{x}, y) = 1) \leftrightarrow f_\beta(\vec{x}, g_{\beta,t}(\vec{x})) = 1)$ .*

*Proof.* We will show separately that:

$$\begin{aligned} U_{PV}^i \vdash \forall \vec{x} ((\exists y \leq t(\vec{x}) \ f_\beta(\vec{x}, y) = 1) \rightarrow f_\beta(\vec{x}, g_{\beta,t}(\vec{x})) = 1), & \quad (\text{Case 1}) \\ U_{PV}^i \vdash \forall \vec{x} (f_\beta(\vec{x}, g_{\beta,t}(\vec{x})) = 1 \rightarrow (\exists y \leq t(\vec{x}) \ f_\beta(\vec{x}, y) = 1)). & \quad (\text{Case 2}) \end{aligned}$$

- **Case 1:** It's easy to see that the sentence we want to prove (in  $U_{PV}^i$ ) is logically equivalent to the following universal sentence:  $\forall \vec{x} \forall y \leq t(\vec{x}) (f_\beta(\vec{x}, y) = 1 \rightarrow f_\beta(\vec{x}, g_{\beta,t}(\vec{x})) = 1)$  (\*). Furthermore, (\*) is a true universal sentence in the standard model by the definition of  $g_{\beta,t}^{\mathbb{N}}$  and  $f_\beta^{\mathbb{N}}$ . Therefore  $U_{PV}^i$  proves (\*).
- **Case 2:** By the definition of  $g_{\beta,t}^{\mathbb{N}}$ , the universal sentence  $\forall \vec{x} \ g_{\beta,t}(\vec{x}) \leq t(\vec{x})$  is true in the standard model, which further means that  $U_{PV}^i \vdash \forall \vec{x} \ g_{\beta,t}(\vec{x}) \leq t(\vec{x})$ . This sentence logically implies the sentence we want to prove in  $U_{PV}^i$ .  $\square$

**Lemma 3.17** (Defining Axioms of  $f_\alpha$ ). *For every  $i \geq 1$  and  $(\Pi_{i-1}^b \cup \Sigma_{i-1}^b)$ -formula  $\alpha(\vec{z})$  in the language  $\mathcal{L}_{PV}$ ,  $U_{PV}^i \vdash \forall \vec{z} (\alpha(\vec{z}) \leftrightarrow f_\alpha(\vec{z}) = 1)$ .*

*Proof.* Fix any  $i \geq 1$ . Let  $\varphi_\alpha \triangleq \forall \vec{z} (\alpha(\vec{z}) \leftrightarrow f_\alpha(\vec{z}) = 1)$ . We firstly prove that  $U_{PV}^i \vdash \varphi_\alpha$  for every bounded  $\mathcal{L}_{PV}$ -formula  $\alpha(\vec{z}) = Q_1x_1 \leq t_1(\vec{z}) \ Q_2x_2 \leq t_2(\vec{z}) \ \dots \ Q_kx_k \leq t_k(\vec{z}) \ \phi(\vec{z}, x_1, \dots, x_k)$ , where  $\phi$  is quantifier free,  $k \leq i-1$ ,  $Q_i \in \{\forall, \exists\}$ , and  $Q_i \neq Q_{i+1}$  for every  $i \in [k-1]$ . We will prove this by induction over  $k$ .

- **(Case 0).** Assume that  $k = 0$  and  $\alpha(\vec{z})$  is a quantifier-free formula. Then  $\varphi_\alpha$  is a universal sentence. Furthermore, by the definition of the interpretation of  $f_\alpha$  over the standard model, we know that  $\mathbb{N} \models \varphi_\alpha$ , which means that  $U_{PV}^i \vdash \varphi_\alpha$ .
- **(Case 1).** Assume that  $\alpha(\vec{z}) = \forall x \leq t(\vec{z}) \alpha'(x, \vec{z})$ . In such case,  $i \geq 2$ . By the induction hypothesis,  $U_{PV}^i \vdash \varphi_{\alpha'}$ . To show that  $U_{PV}^i \vdash \varphi_\alpha$  it is sufficient to prove that  $U_{PV}^i \vdash \forall \vec{z} (f_\alpha(\vec{z}) = 1 \rightarrow \alpha(\vec{z}))$  and  $U_{PV}^i \vdash \forall \vec{z} (\alpha(\vec{z}) \rightarrow f_\alpha(\vec{z}) = 1)$ . Now we prove them separately.

(i) Since  $U_{PV}^i \vdash \varphi_{\alpha'}$ , we know that  $U_{PV}^i$  proves  $\forall \vec{z} \forall x \leq t(\vec{z}) (f_{\alpha'}(x, \vec{z}) = 1 \rightarrow \alpha'(x, \vec{z}))$  ( $\star$ ). Consider the universal sentence  $\psi \triangleq \forall \vec{z} \forall x \leq t(\vec{z}) (f_\alpha(\vec{z}) = 1 \rightarrow f_{\alpha'}(x, \vec{z}) = 1)$ . By the definition of the interpretations of  $f_\alpha$  and  $f_{\alpha'}$ ,  $\psi$  is a true sentence, therefore  $U_{PV}^i \vdash \psi$  ( $\diamond$ ). Combining ( $\star$ ) and ( $\diamond$ ) we get that

$$U_{PV}^i \vdash \forall \vec{z} \forall x \leq t(\vec{z}) (f_\alpha(\vec{z}) = 1 \rightarrow \alpha'(x, \vec{z})).$$

This means that  $U_{PV}^i \vdash \forall \vec{z} (f_\alpha(\vec{z}) = 1 \rightarrow \alpha(\vec{z}))$ .

(ii) Recall that we need to show that  $U_{PV}^i \vdash \forall \vec{z} ((\forall x \leq t(\vec{z}) \alpha'(x, \vec{z})) \rightarrow f_\alpha(\vec{z}) = 1)$ . Since  $U_{PV}^i \vdash \varphi_{\alpha'}$ , it is sufficient to prove that

$$U_{PV}^i \vdash \forall \vec{z} ((\forall x \leq t(\vec{z}) f_{\alpha'}(\vec{z}, x) = 1) \rightarrow f_\alpha(\vec{z}) = 1).$$

Since  $\alpha$  is a  $\Pi_{i-1}^b$ -formula of the form above,  $\neg \alpha'$  is a  $\Sigma_{i-2}^b \cup \Pi_{i-2}^b$ -formula. By Lemma 3.15,  $U_{PV}^i \vdash f_{\alpha'}(\vec{z}, x) = 1 \leftrightarrow f_{\neg \alpha'}(\vec{z}, x) = 0$ . So we only need to prove that

$$U_{PV}^i \vdash \forall \vec{z} ((\forall x \leq t(\vec{z}) f_{\neg \alpha'}(\vec{z}, x) = 0) \rightarrow f_\alpha(\vec{z}) = 1).$$

By Lemma 3.16, we know that  $U_{PV}^i \vdash (\exists x \leq t(\vec{z}) f_{\neg \alpha'}(\vec{z}, x) = 1) \leftrightarrow f_{\neg \alpha'}(\vec{z}, g_{\neg \alpha', t}(\vec{z})) = 1$ , which means we only need to prove that

$$U_{PV}^i \vdash \forall \vec{z} (f_{\neg \alpha'}(\vec{z}, g_{\neg \alpha', t}(\vec{z})) = 0 \rightarrow f_\alpha(\vec{z}) = 1). \quad (3)$$

By considering the interpretations of  $f_\alpha$ ,  $f_{\neg \alpha'}$ , and  $g_{\neg \alpha', t}$  in the standard model, it follows that the universal sentence (3) is true in the standard model. Therefore it is provable in  $U_{PV}^i$ . This completes this case.

- **(Case 2).** Assume that  $\alpha(\vec{z}) = \exists x \leq t(\vec{z}) \alpha'(x, \vec{z})$ . Let  $\bar{\alpha}(\vec{z})$  be the formula obtained by pushing the negation in  $\neg \alpha(\vec{z})$  into the quantifiers. Note that  $\vdash \neg \alpha(\vec{z}) \leftrightarrow \bar{\alpha}(\vec{z})$ . By applying Case 1, we can show that

$$U_{PV}^i \vdash \forall \vec{z} (f_{\bar{\alpha}}(\vec{z}) = 1 \leftrightarrow \neg \alpha(\vec{z})).$$

Since  $\forall \vec{z} (f_{\bar{\alpha}}(\vec{z}) = 1 \leftrightarrow f_\alpha(\vec{z}) \neq 1)$  is a universal sentence that is true in the standard model, we know that it is provable in  $U_{PV}^i$ , which further implies that

$$U_{PV}^i \vdash \forall \vec{z} (f_\alpha(\vec{z}) \neq 1 \leftrightarrow \neg \alpha(\vec{z})).$$

This yields  $U_{PV}^i \vdash \varphi_\alpha$ .

Now we consider the case when  $\alpha$  is an arbitrary  $(\Pi_{i-1}^b \cup \Sigma_{i-1}^b)$ -formula. By Lemma 3.14, we can see that  $U_{PV}^i \vdash \forall \vec{z} (\alpha(\vec{z}) \leftrightarrow \alpha^{\text{norm}}(\vec{z}))$ . According the discussion above, we know that  $U_{PV}^i \vdash \forall \vec{z} (\alpha^{\text{norm}}(\vec{z}) \leftrightarrow f_{\alpha^{\text{norm}}}(\vec{z}) = 1)$ . Moreover, we have that  $U_{PV}^i \vdash \forall \vec{z} (f_\alpha(\vec{z}) = 1 \leftrightarrow f_{\alpha^{\text{norm}}}(\vec{z}) = 1)$ , since this is a true universal sentence in the standard model. It follows from the provability of these three sentences that  $U_{PV}^i \vdash \varphi_\alpha$ , as desired.  $\square$

**Theorem 3.18.** *For every  $i \geq 1$  and  $\mathcal{L}_{\text{PV}}$ -sentence  $\varphi$ , if  $T_{\text{PV}}^i \vdash \varphi$ , then  $U_{\text{PV}}^i \vdash \varphi$ .*

*Proof.* To prove this lemma, it is sufficient to show that for every  $\varphi \in T_{\text{PV}}^i$ ,  $U_{\text{PV}}^i \vdash \varphi$ . Let  $\varphi = \forall \vec{x} \alpha(\vec{x})$  be an axiom of  $T_{\text{PV}}^i$ , where  $\alpha(\vec{x})$  is a  $\Sigma_{i-1}^b$ -formula. By Lemma 3.17, we only need to show that  $U_{\text{PV}}^i$  proves  $\forall \vec{x} f_\alpha(\vec{x}) = 1$ . This follows directly from the fact that  $\forall \vec{x} f_\alpha(\vec{x}) = 1$  is a true universal sentence in the standard model.  $\square$

**Complexity of the function symbols in  $\mathcal{L}_{\text{PV}}^i$ .** As we discussed in Section 1.2, we will extract a KPT-style student-teacher game from the provability of the lower bound sentence in the universal theory  $U_{\text{PV}}^i$ . In this step, the complexity of the student is determined by the complexity of the standard interpretations of the function symbols in the language  $\mathcal{L}_{\text{PV}}^i$ , which consists of both the polynomial-time computable functions (i.e. the symbols in  $\mathcal{L}_{\text{PV}}$ ) and the new function symbols  $f_\alpha$  and  $g_{\beta,t}$ . Now we determine the complexity of the functions  $f_\alpha$  and Skolem functions  $g_{\beta,t}$ .

**Lemma 3.19.** *Let  $i \geq 1$ . For every function symbol  $f_\alpha$  in  $\mathcal{L}_{\text{PV}}^i$ ,  $f_\alpha^\mathbb{N} : \mathbb{N} \rightarrow \{0, 1\}$  is the characteristic function of a language in  $\Pi_{i-1}^p \cup \Sigma_{i-1}^p$ .*

*Proof.* Recall that each  $f_\alpha$  is introduced for a  $(\Pi_{i-1}^b \cup \Sigma_{i-1}^b)$ -formula  $\alpha(\vec{z})$  with language  $\mathcal{L}_{\text{PV}}$  such that  $f_\alpha^\mathbb{N}$  is the characteristic function of  $\alpha^\mathbb{N}$ , i.e., for every  $\vec{m} \in \vec{\mathbb{N}}$ ,  $f_\alpha^\mathbb{N}(\vec{m}) = 1$  if and only if  $\alpha^\mathbb{N}(\vec{m})$  holds. Since  $\alpha$  is a bounded formula and the initial function symbols and relation symbols, when interpreted in the standard model, are polynomial-time computable, it is not hard to see that  $\alpha^\mathbb{N} \in \Pi_{i-1}^p \cup \Sigma_{i-1}^p$ .  $\square$

**Lemma 3.20.** *Let  $i \geq 2$ . For every function symbol  $g_{\beta,t}$  in  $\mathcal{L}_{\text{PV}}^i$ ,  $g_{\beta,t}^\mathbb{N} \in \text{FP}^{\Sigma_{i-1}^p}$ .*

*Proof.* Recall that  $g_{\beta,t}$  is introduced for every  $\Sigma_{i-1}^b$ -formula  $\beta$  and term  $t$  in the language  $\mathcal{L}_{\text{PV}}$ , and that  $g_{\beta,t}^\mathbb{N}(\vec{x})$  finds the minimum  $y^*$  such that  $\beta^\mathbb{N}(\vec{x}, y^*)$  holds if there is  $y \leq t(\vec{x})$  such that  $\beta^\mathbb{N}(\vec{x}, y)$  holds, or outputs 0 otherwise. Note that using a  $\Sigma_{i-1}^p$  oracle we can decide for  $0 \leq l \leq r \leq t(\vec{x})$  whether there exists  $y \in [l, r]$  such that  $\beta^\mathbb{N}(\vec{x}, y)$  holds. So we can perform a binary search over  $[0, t(\vec{x})]$  to find the minimum  $y^*$  such that  $\beta^\mathbb{N}(\vec{x}, y^*)$  holds or detect that no such element exists. This is an  $\text{FP}^{\Sigma_{i-1}^p}$  computation for every  $i \geq 2$ .  $\square$

**Theorem 3.21.** *Let  $i \geq 1$ . For every  $\mathcal{L}_{\text{PV}}^i$ -term  $t(x_1, \dots, x_\ell)$ , we have  $t^\mathbb{N}(x_1, \dots, x_\ell) \in \text{FP}^{\Sigma_{i-1}^p}$ .*

*Proof.* This directly follows from Lemma 3.19 and Lemma 3.20.  $\square$

Theory  $U_{\text{PV}}^i$  has almost all properties needed for the proof of our results, except that it is not necessarily closed under if-then-else (Definition 2.2). This is desirable as it simplifies the statement of our witnessing result and its proof. For this reason, we further modify  $U_{\text{PV}}^i$  to guarantee this property.

**Theory  $UT_{\text{PV}}^i$  and Language  $\mathcal{L}_{\text{UT}}^i$ .** Let  $i \geq 1$ , and consider the language  $\mathcal{L}_{\text{PV}}^i$  introduced before. We extend  $\mathcal{L}_{\text{PV}}^i$  as follows. For every  $k \geq 1$  and function  $f : \mathbb{N}^k \rightarrow \mathbb{N}$  in  $\text{FP}^{\Sigma_{i-1}^p}$ , we introduce a new function symbol  $f_{\text{UT}}$ . Then, we let

$$\mathcal{L}_{\text{UT}}^i \triangleq \mathcal{L}_{\text{PV}}^i \cup \{f_{\text{UT}} \mid f \in \text{FP}^{\Sigma_{i-1}^p}\}.$$

Given  $\mathcal{L}_{\text{UT}}^i$ , we define  $UT_{\text{PV}}^i$  as the theory of all universal sentences in  $\mathcal{L}_{\text{UT}}^i$  that are true in the standard model.

**Theorem 3.22 (Main Properties of  $UT_{\text{PV}}^i$ ).** *For every  $i \geq 1$ , the theory  $UT_{\text{PV}}^i$  satisfies the following properties:*

- (i)  $\text{UT}_{\text{PV}}^i$  is a universal theory.
- (ii) Every  $\mathcal{L}_{\text{PV}}^i$ -sentence provable in  $\text{U}_{\text{PV}}^i$  is also provable in  $\text{UT}_{\text{PV}}^i$ .
- (iii) Every  $\mathcal{L}_{\text{PV}}$ -sentence provable in  $\text{T}_{\text{PV}}^i$  is also provable in  $\text{UT}_{\text{PV}}^i$ .
- (iv) Let  $t$  be an arbitrary  $\mathcal{L}_{\text{UT}}^i$ -term, and consider its interpretation  $t^{\mathbb{N}}: \mathbb{N}^k \rightarrow \mathbb{N}$  over the standard model. Then  $t^{\mathbb{N}} \in \text{FP}^{\Sigma_{i-1}^p}$ .
- (v)  $\text{UT}_{\text{PV}}^i$  is closed under if-then-else.
- (vi)  $\text{UT}_{\text{PV}}^i$  is sound, i.e., every sentence provable in  $\text{UT}_{\text{PV}}^i$  is true over  $\mathbb{N}$ .

*Proof.* We prove each item in turn.

- (i) This is immediate from the definition of the theory.
- (ii) Let  $\varphi$  be an  $\mathcal{L}_{\text{PV}}^i$ -sentence provable in  $\text{U}_{\text{PV}}^i$ . It is enough to argue that every axiom of  $\text{U}_{\text{PV}}^i$  is provable in  $\text{UT}_{\text{PV}}^i$ . Since  $\text{U}_{\text{PV}}^i$  is the theory consisting of all universal true sentences (over the standard model) in  $\mathcal{L}_{\text{PV}}$ ,  $\mathcal{L}_{\text{PV}}^i \subseteq \mathcal{L}_{\text{UT}}^i$ , and  $\text{UT}_{\text{PV}}^i$  is the theory of all universal sentences in  $\mathcal{L}_{\text{UT}}^i$  that are true in the standard model, the result is immediate.
- (iii) Let  $\varphi$  be an  $\mathcal{L}_{\text{PV}}$ -sentence provable in  $\text{T}_{\text{PV}}^i$ . It follows from Theorem 3.18 that  $\varphi$  is provable in  $\text{U}_{\text{PV}}^i$ . Consequently, the claim follows from the previous item.
- (iv) This follows from Theorem 3.21, the definition of  $\mathcal{L}_{\text{UT}}^i$ , and the closure of the functions in  $\text{FP}^{\Sigma_{i-1}^p}$  under composition.
- (v) To show this, let  $\varphi(x_1, \dots, x_k)$  be a quantifier-free  $\mathcal{L}_{\text{UT}}^i$ -formula, and consider  $\mathcal{L}_{\text{UT}}^i$ -terms  $t_1(x_1, \dots, x_k)$  and  $t_2(x_1, \dots, x_k)$ . We must prove that there exists an  $\mathcal{L}_{\text{UT}}^i$ -term  $t(x_1, \dots, x_k)$  such that

$$\text{UT}_{\text{PV}}^i \vdash (t(\vec{x}) = t_1(\vec{x}) \wedge \varphi(\vec{x})) \vee (t(\vec{x}) = t_2(\vec{x}) \wedge \neg\varphi(\vec{x})). \quad (4)$$

Consider the interpretations of terms  $t_1^{\mathbb{N}}, t_2^{\mathbb{N}}: \mathbb{N}^k \rightarrow \mathbb{N}$  over the standard model. Let  $f: \mathbb{N}^k \rightarrow \mathbb{N}$  be the function defined as follows:

$$f(\vec{a}) = \begin{cases} t_1^{\mathbb{N}}(\vec{a}) & \text{if } \varphi^{\mathbb{N}}(\vec{a}) \text{ is true;} \\ t_2^{\mathbb{N}}(\vec{a}) & \text{otherwise.} \end{cases}$$

Since  $\varphi$  is a quantifier-free formula, thanks to Item (iii), it is easy to see that  $f \in \text{FP}^{\Sigma_{i-1}^p}$ . Consequently, the corresponding function symbol  $f_{\text{UT}} \in \mathcal{L}_{\text{UT}}^i$ . Take  $t$  as  $f_{\text{UT}}$ . It follows from the definition of  $f$  and of  $t$  that, for every  $\vec{a} \in \mathbb{N}^k$ ,

$$\mathbb{N} \models (t(\vec{a}) = t_1(\vec{a}) \wedge \varphi(\vec{a})) \vee (t(\vec{a}) = t_2(\vec{a}) \wedge \neg\varphi(\vec{a})).$$

Since the formula above is free of quantifiers, the definition of  $\text{UT}_{\text{PV}}^i$  immediately yields Equation (4).

- (vi) This is obvious from its definition.  $\square$

## 4 Witnessing Theorems for General Formulas

In this section, we establish a witnessing theorem that works for sentences of arbitrarily high quantifier complexity. As explained in Section 1, the result is a key step in our proof that strong complexity lower bounds cannot be established in  $\text{T}_{\text{PV}}^i$ . While it is possible to obtain a general witnessing result that holds for an arbitrary universal theory, due to our main applications we restrict our attention to theories of bounded arithmetic.



## 4.1 A game-theoretic witnessing theorem

Let  $\mathcal{T}$  be a universal bounded theory over the vocabulary  $\mathcal{L}$ . Let  $\varphi(x)$  be a bounded  $\mathcal{L}$ -formula defined as

$$\begin{aligned} \varphi(x) \triangleq & \exists y_1 \leq t_1(x) \forall x_1 \leq s_1(x, y_1) \exists y_2 \leq t_2(x, y_1, x_1) \dots \forall x_{k-1} \leq s_{k-1}(x, y_1, x_1, \dots, y_{k-1}) \\ & \exists y_k \leq t_k(x, y_1, x_1, \dots, y_{k-1}, x_{k-1}) \forall x_k \leq s_k(x, y_1, x_1, \dots, y_k) \phi(x, x_1, \dots, x_k, y_1, \dots, y_k), \end{aligned}$$

where  $\phi(x, \vec{x}, \vec{y})$  is a quantifier-free  $\mathcal{L}$ -formula.

**The Evaluation Game.** We consider an interactive game between two players, the *truthifier* (associated with existential quantifiers in  $\varphi$ ) and the *falsifier* (associated with universal quantifiers in  $\varphi$ ). A *board* is defined as a pair  $(\mathcal{M}, n_0)$ , where  $\mathcal{M}$  is a structure over  $\mathcal{L}$  such that  $\mathcal{M} \models \mathcal{T}$ , and  $n_0 \in \mathcal{M}$  is an element of its domain.<sup>16</sup> The *evaluation game* for the formula  $\varphi(x)$  on the board  $(\mathcal{M}, n_0)$  is played as follows: in the  $i$ -th round of the game ( $1 \leq i \leq k$ ), the truthifier firstly chooses an assignment  $m_i \in \mathcal{M}$  for  $y_i$  such that  $m_i \leq t_i(n_0, m_1, n_1, \dots, m_{i-1}, n_{i-1})$ , then the falsifier chooses an assignment  $n_i \in \mathcal{M}$  for  $x_i$  such that  $n_i \leq s_i(n_0, m_1, n_1, \dots, m_i)$ . The truthifier *wins* if and only if  $\phi(x/n_0, \vec{x}/\vec{n}, \vec{y}/\vec{m})$  holds in  $\mathcal{M}$ .

The *transcript* of a game given strategies  $\tau^t$  for the truthifier and  $\tau^f$  for the falsifier, denoted by  $\langle \tau^t : \tau^f \rangle$ , is a pair  $(\vec{n}, \vec{m})$  of sequences that records the moves for both players.<sup>17</sup> A *partial transcript* is a prefix of a transcript. A partial transcript is *valid* if all elements  $m_i$  and  $n_i$  respect the corresponding upper bounds (in  $\mathcal{M}$ ) given by functions  $t_i$  and  $s_i$ . A strategy  $\tau^t$  for the truthifier is said to *beat* a strategy  $\tau^f$  for the falsifier (w.r.t. a given board and formula) if the truthifier wins in the transcript  $\langle \tau^t : \tau^f \rangle$ . A *strategy* for a player is defined in the natural way, i.e., as a function that produces the next assignment given a partial transcript of the game. Equivalently, since we will consider games with only a fixed number of rounds, one can describe a strategy as a finite sequence of functions of the form  $f : \mathcal{M}^i \rightarrow \mathcal{M}$ , for appropriate values  $i \leq 2k$ .

We will consider games that are played in a more general setting. Roughly speaking, we allow the truthifier and falsifier to simultaneously play different evaluation games over the same board  $(\mathcal{M}, n_0)$ . The truthifier has a positional advantage over the falsifier: it can decide where to make the next move, i.e., by either making the next move in one of the current games or starting a new game over the board  $(\mathcal{M}, n_0)$  or playing differently some earlier play, which creates a new game from there but maintains the existing game plays. The falsifier must respond to that move in the corresponding game. *Crucially, the next assignment selected by each player now depends on previous plays in all games.* The formal details are provided next.

**The Tree Exploration Game.** A *partial game tree*  $T = (V, E, \gamma)$  (where  $(V, E)$  is a directed rooted tree and  $\gamma : E \rightarrow \mathcal{M} \times \mathcal{M}$ ) of the evaluation game for  $\varphi$  on the board  $(\mathcal{M}, n_0)$  is defined as a finite rooted tree where each edge  $e \in E(T)$  is labeled with a pair  $(m, n)$  of elements of  $\mathcal{M}$  and such that, for every node  $u \in V(T)$ , the concatenation of each pair of elements labelling the edges on the root-to- $u$  path is a prefix of a valid transcript of the evaluation game of  $\varphi(x)$  on the board  $(\mathcal{M}, n_0)$ . More precisely, if the pairs labelling the edges from the root to  $u$  are  $(m_1, n_1), (m_2, n_2), \dots, (m_i, n_i)$ , then  $(m_1, n_1, m_2, n_2, \dots, m_i, n_i)$  is a valid partial transcript of the evaluation game, i.e., for all  $j \in [i]$ ,  $\mathcal{M} \models m_j \leq t_j(n_0, m_1, n_1, \dots, m_{j-1}, n_{j-1})$  and  $\mathcal{M} \models n_j \leq s_j(n_0, m_1, n_1, \dots, m_j)$ . Note that if  $\mathcal{M}$  is the standard model then a partial game tree of the evaluation game is a finite upper part of the (exponential size) complete game tree of the evaluation game.

<sup>16</sup>For a concrete example, think of  $\mathcal{M} = (\mathbb{N}, \leq^{\mathbb{N}}, +^{\mathbb{N}}, \times^{\mathbb{N}}, \dots)$ .

<sup>17</sup>For convenience, we might also write the transcript as  $(m_1, n_1, \dots, m_k, n_k)$ . The moves of each player will always be clear in each context.

Let  $T$  be a partial game tree of  $\varphi$  and  $(\mathcal{M}, n_0)$  be a board. The *tree exploration game* starting from  $T$  on  $(\mathcal{M}, n_0)$  is played as follows. In each *round*, first the truthifier chooses a node  $u$  from  $T$  (not necessarily a leaf) and an element  $m \in \mathcal{M}$ , then the falsifier chooses an element  $n \in \mathcal{M}$ . This creates a child of  $u$  and a corresponding directed edge labeled by  $(m, n)$ . Note that when playing each round of the tree exploration game both players should guarantee that the new partial game tree is always a valid partial game tree, i.e.,  $m$  and  $n$  should satisfy the corresponding inequalities. The *size* of a partial game tree  $T$  is given by  $|T(V)|$ .

The truthifier *wins* the tree exploration game if there is a node in the current partial game tree that is a winning node for the truthifier, that is, the concatenation of the pairs of elements labelling the edges on the root-to- $u$  path forms a winning transcript of the truthifier in the evaluation game of  $\varphi(x)$  on the board  $(\mathcal{M}, n_0)$ . The *tree exploration game of  $\varphi(x)$*  is defined as the tree exploration game starting from a partial game tree containing only the root node. We refer to Figure 1 for an example of the tree exploration game.

Recall that  $\mathcal{L}$  is the vocabulary of the universal (bounded) theory  $\mathcal{T}$ . The main result established in this section shows the existence of a “computationally bounded” winning strategy for the truthifier from a  $\mathcal{T}$ -proof of  $\varphi$ , i.e., the strategy can be computed by  $\mathcal{L}$ -terms. In addition, the strategy is universal, in the sense that it is specified by  $\mathcal{L}$ -terms that are independent of the board  $(\mathcal{M}, n_0)$ . Finally, the location of each play of the truthifier in the partial game tree is fixed in advance and does not depend on the strategy of the falsifier nor on the board  $(\mathcal{M}, n_0)$ . (The elements selected by the truthifier depend on the previous plays of the truthifier and falsifier.) This means that the trees in the sequence of partial game trees are fixed in advance.

**$\mathcal{L}$ -Strategies for the Tree Exploration Game.** An  $\mathcal{L}$ -*quasi-strategy* of the truthifier of *length*  $\ell \in \mathbb{N}$  and initial tree size  $d$  is a sequence  $\tau = \langle p_1, r_1, p_2, r_2, \dots, p_\ell, r_\ell \rangle$ , where each  $p_i$  is an  $\mathcal{L}$ -term and each  $r_i \in \mathbb{N}$  is such that  $1 \leq r_i < d + i$ . Let  $(\mathcal{M}, n_0)$  be a board and  $T$  be a partial game tree on this board with  $V(T) = \{1, 2, \dots, d\}$ . The strategy for the tree exploration game starting from the partial game tree  $T$  induced by  $\tau$  proceeds as follows:

- In the  $i$ -th move, the truthifier introduces a node numbered  $d + i$  as a child of the node  $r_i$  and chooses the element  $v_i \triangleq p_i^{\mathcal{M}}(n_0, \Gamma) \in \mathcal{M}$ , where  $\Gamma$  is the sequence of  $\mathcal{M}$ -elements chosen by the players in previous rounds (i.e.,  $v_1, \dots, v_{i-1}$  and the falsifier’s moves  $w_1, \dots, w_{i-1}$ ).

Note that an arbitrary  $\mathcal{L}$ -*quasi-strategy* might induce an invalid move  $v_i = p_i^{\mathcal{M}}(n_0, \Gamma)$  that violates the desired upper bound on  $v_i$ , depending on the moves of the falsifier. We say that an  $\mathcal{L}$ -*quasi-strategy* of the truthifier is an  $\mathcal{L}$ -*strategy* if for every board  $(\mathcal{M}, n_0)$  the resulting partial game trees are valid for every valid strategy of the falsifier.

Finally, a length- $\ell$   $\mathcal{L}$ -strategy is said to be a *universal winning strategy* if the truthifier wins within  $\ell$  moves against all valid strategies (not necessarily  $\mathcal{L}$ -strategies) of the falsifier on any board  $(\mathcal{M}, n_0)$ . Note that the falsifier’s strategy is a function of the board  $(\mathcal{M}, n_0)$ , partial game tree  $T = (V, E, \gamma)$  (which includes all moves from previous rounds), and the move of the truthifier in the current round.

**Theorem 4.1** (Game-Theoretic Witnessing Theorem). *Let  $\mathcal{T}$  be a universal bounded theory with vocabulary  $\mathcal{L}$  that is closed under if-then-else. Let  $\varphi$  be a bounded  $\mathcal{L}$ -formula of the form*

$$\begin{aligned} \varphi(x) \triangleq & \exists y_1 \leq t_1(x) \forall x_1 \leq s_1(x, y_1) \exists y_2 \leq t_2(x, y_1, x_1) \dots \forall x_{k-1} \leq s_{k-1}(x, y_1, x_1, \dots, y_{k-1}) \\ & \exists y_k \leq t_k(x, y_1, x_1, \dots, y_{k-1}, x_{k-1}) \forall x_k \leq s_k(x, y_1, x_1, \dots, y_k) \phi(x, x_1, \dots, x_k, y_1, \dots, y_k), \end{aligned}$$

where  $\phi(x, \vec{x}, \vec{y})$  is a quantifier-free  $\mathcal{L}$ -formula. Then  $\mathcal{T} \vdash \forall x \varphi(x)$  if and only if there is a universal winning  $\mathcal{L}$ -strategy of length  $O(1)$  for the truthifier in the corresponding tree exploration game of  $\varphi(x)$ .

## 4.2 A cut-free sequent calculus

We assume familiarity with basic aspects of proof theory and Gentzen-style sequent calculi for first-order logic. We refer the reader to the standard textbook [TS00] for the necessary background.

Our main technical tool to prove the game-theoretic witnessing theorem is a cut-free sequent calculus for classical first-order logic known as G3c (see Sections 3 and 4 of [TS00]).<sup>18</sup> We follow the notation and conventions from [TS00], unless explicitly stated otherwise.

For simplicity, we will work with the (complete) set  $\{\perp, \rightarrow\}$  of connectives. Other Boolean connectives can be expressed using these in the usual way. For this set of connectives, the system G3c contains the following rules (see [TS00, Section 3.5]).<sup>19</sup>

### Sequent Calculus G3c:

$$\begin{array}{ll}
 \text{Ax: } \frac{}{\Gamma, P \Rightarrow \Delta, P} & \text{L}\perp: \frac{}{\Gamma, \perp \Rightarrow \Delta} \\
 \text{L}\rightarrow: \frac{\Gamma \Rightarrow \Delta, \alpha \quad \Gamma, \beta \Rightarrow \Delta}{\Gamma, \alpha \rightarrow \beta \Rightarrow \Delta} & \text{R}\rightarrow: \frac{\Gamma, \alpha \Rightarrow \Delta, \beta}{\Gamma \Rightarrow \Delta, \alpha \rightarrow \beta} \\
 \text{L}\forall: \frac{\Gamma, \alpha(x/t), \forall x \alpha(x) \Rightarrow \Delta}{\Gamma, \forall x \alpha(x) \Rightarrow \Delta} & \text{R}\forall: \frac{\Gamma \Rightarrow \Delta, \alpha(x/y)}{\Gamma \Rightarrow \Delta, \forall x \alpha(x)} \\
 \text{L}\exists: \frac{\Gamma, \alpha(x/y) \Rightarrow \Delta}{\Gamma, \exists x \alpha(x) \Rightarrow \Delta} & \text{R}\exists: \frac{\Gamma \Rightarrow \Delta, \alpha(x/t), \exists x \alpha(x)}{\Gamma \Rightarrow \Delta, \exists x \alpha(x)}
 \end{array}$$

**Notes.** Recall that in a sequent  $\Sigma \Rightarrow \Pi$  both  $\Sigma$  and  $\Pi$  are multisets of formulas. As usual,  $\Sigma$  is called the *antecedent* and  $\Pi$  is called the *succedent*. In the rules above,  $\alpha$  and  $\beta$  are formulas. In the rule (Ax),  $P$  is an atomic predicate, i.e.,  $P$  is of the form  $R(t_1, \dots, t_k)$ , where  $R$  is a relation symbol and  $t_1, \dots, t_k$  are arbitrary terms. The variable  $y$  appearing in rules (R $\forall$ ) and (L $\exists$ ) has no free occurrences in the conclusion, i.e., in the formulas in  $\Gamma \cup \Delta \cup \{\forall x \alpha\}$  or  $\Gamma \cup \Delta \cup \{\exists x \alpha\}$ . The formula denoted by  $\alpha(x)$  can contain free variables besides  $x$ . A proof tree of a sequent  $\Sigma \Rightarrow \Pi$  is a tree of the applications of the eight rules, where the conclusion (root of the tree) is  $\Sigma \Rightarrow \Pi$  and the leaves have no premise (i.e., they are applications of (Ax) or (L $\perp$ )).

**Naming convention.** Following the naming convention of [TS00], for all the eight rules, we call the formula in the conclusion besides  $\Gamma$  and  $\Delta$  the *principal formula*. The formulas in  $\Gamma$  and  $\Delta$  in both the premise and the conclusion are called *side formulas*. The formulas in  $\Gamma, \Delta$  in the conclusion of (Ax) and (L $\perp$ ) are called *weak formulas*. The formulas in the premise that are neither side nor weak are called *active formulas*; in particular, the formulas in the premises that are identical to the principal formulas (i.e.,  $\forall x \alpha(x)$  in (L $\forall$ ) and  $\exists x \alpha(x)$  in (R $\exists$ )) are called *Kleene copies*. The variable  $y$  in (R $\forall$ ) or (L $\exists$ ) is called the *proper variable* of  $\alpha$ .

<sup>18</sup>To give some context, in [TS00] the system G1c is close to the original sequent calculus LK of Gentzen, G2c is a variant with weakening absorbed into the logical rules, and G3c has both weakening and contraction absorbed into the rules and axioms.

<sup>19</sup>Restricting the rules of G3c to our set of connectives can be done without loss of generality by the “separation property” of G3c (see the first corollary in [TS00, Section 4.2]).

**Bounded variable renaming.** As in [TS00], we identify the formulas that differ only by a renaming of the bounded variables in the rest of the section. Strictly speaking, this means that formulas appeared in the sequents are actually representing an equivalence class of the formulas up to a bounded variable renaming. Under this convention, we can have the completeness of G3c:

**Theorem 4.2** (Completeness of G3c [TS00, Chapter 4]). *Let  $\Gamma, \Delta$  be finite multiset of formulas. If  $\bigwedge \Gamma \rightarrow \bigvee \Delta$  is a tautology, then  $\Gamma \Rightarrow \Delta$  is a provable sequent.*

**Dealing with Non-logical Axioms.** The sequent calculus G3c defined above is for first-order logic without non-logical axioms and the axioms for equality. To deal with the theory  $T_{PV}^i$  that contains (infinitely many) non-logical axioms, we can simply use the compactness theorem of first-order logic.

**Theorem 4.3** (Compactness Theorem for First-Order Logic). *Let  $\Gamma$  be a set of formulas, and  $\varphi$  be a formula. Then  $\Gamma \vdash \varphi$  if and only if there is a finite multiset  $\Gamma_0$  of formulas in  $\Gamma$  such that  $\Gamma_0 \vdash \varphi$ .*

The compactness theorem and the completeness of G3c imply that for every theory  $\mathcal{T}$ , if  $\mathcal{T} \vdash \alpha$ , then there exists a finite set  $\Gamma$  of axioms in  $\mathcal{T}$  such that  $\Gamma \Rightarrow \alpha$  is a provable sequent. In particular, when  $\mathcal{T}$  is a universal theory, as in the statement of Theorem 4.1, we can assume w.l.o.g. that the sentences in  $\Gamma$  are explicitly written in the form  $\forall \vec{x} \varphi(\vec{x})$  for some quantifier-free formula  $\varphi$ . We will then extract a universal winning  $\mathcal{L}$ -strategy from the proof tree of the sequent  $\Gamma \Rightarrow \alpha$ .

**Comparison with proof systems in previous works.** Many proof-theoretic results in bounded arithmetic utilize some systems of sequent calculus that contain restricted cut rules (see, e.g., the free-cut-free system in [Bus86] and anchored LK in [CN10]). Those systems are designed to treat the theories as a part of the proof system, while the system G3c we use do not have this feature.

**Representation of bounded formulas.** To simplify the structure of the proof trees, we need to choose a canonical representation of bounded formulas. Let  $\varphi(x)$  be a  $\Sigma_k^b$ -formula in prenex normal form (with bounded quantifiers). We define the following translation  $[\cdot]_{\text{imp}}$  that transforms a bounded formula in prenex normal form into a logically equivalent formula with only unbounded quantifiers:

- If  $\varphi(\vec{x})$  is quantifier free,  $[\varphi(\vec{x})]_{\text{imp}} \triangleq \varphi(\vec{x})$ .
- If  $\varphi(\vec{x}) = \forall y \leq t(\vec{x}) \phi(\vec{x}, y)$  and  $[\phi]_{\text{imp}} = Q_1 z_1 Q_2 z_2 \dots Q_k z_k \alpha(\vec{x}, y, \vec{z})$ , where  $Q_i \in \{\forall, \exists\}$  for  $i \in [k]$  and  $\alpha$  is quantifier-free, then  $[\varphi(\vec{x})]_{\text{imp}} \triangleq \forall y Q_1 z_1 Q_2 z_2 \dots Q_k z_k (y \leq t(\vec{x}) \rightarrow \alpha(\vec{x}, y, \vec{z}))$ .
- If  $\varphi(\vec{x}) = \exists y \leq t(\vec{x}) \phi(\vec{x}, y)$  and  $[\phi]_{\text{imp}} = Q_1 z_1 Q_2 z_2 \dots Q_k z_k \alpha(\vec{x}, y, \vec{z})$ , where  $Q_i \in \{\forall, \exists\}$  for  $i \in [k]$  and  $\alpha$  is quantifier-free, then  $[\varphi(\vec{x})]_{\text{imp}} \triangleq \exists y Q_1 z_1 Q_2 z_2 \dots Q_k z_k (y \leq t(\vec{x}) \wedge \alpha(\vec{x}, y, \vec{z}))$ .<sup>20</sup>

We say a formula  $\varphi$  is *implicitly bounded* if there is a bounded formula  $\psi$  in prenex normal form such that  $\varphi = [\psi]_{\text{imp}}$ .

**Proposition 4.4.** *For every bounded formula  $\varphi(x)$  in prenex normal form (with bounded quantifiers),  $[\varphi(x)]_{\text{imp}}$  is logically equivalent to  $\varphi(x)$ .*

*Proof.* This follows from a simple induction with the logical equivalence  $\alpha \star (Qx \beta) \equiv Qx (\alpha \star \beta)$  for any  $\star \in \{\rightarrow, \wedge\}$  and  $Q \in \{\forall, \exists\}$  when  $x$  has no free occurrence in  $\alpha$ .  $\square$

<sup>20</sup>Note that a conjunction  $\alpha \wedge \beta$  can be written as  $(\alpha \rightarrow (\beta \rightarrow \perp)) \rightarrow \perp$ . As mentioned above, when discussing provability in G3c we will assume w.l.o.g. that the quantifier-free part of  $[\varphi]_{\text{imp}}$  is written using connectives from  $\{\rightarrow, \perp\}$ .

### 4.3 Structural transformations of the proof tree

Let  $\varphi(x) = \exists y_1 \forall x_1 \dots \exists y_k \forall x_k \phi(x, \vec{x}, \vec{y})$  be an implicitly bounded formula with exactly one free variable  $x$  for quantifier-free  $\phi$ . We denote the sub-formula of  $\varphi$  starting from the  $(i + 1)$ -th existential quantifier  $\exists y_{i+1} \forall x_{i+1} \dots \exists y_k \forall x_k \phi(x, \vec{x}, \vec{y})$  as  $\varphi[i]$  (e.g.  $\varphi = \varphi[0]$ ).<sup>21</sup> Let  $\Gamma$  be a finite multiset of universal sentences written explicitly in the form  $\forall \vec{x} \beta(\vec{x})$  for some quantifier-free formula  $\beta(x)$ , i.e., it contains only universal quantifiers and no propositional connective occurs outside of a (universal) quantifier. We will prove some structural lemmas for a proof tree of  $\Gamma \Rightarrow \forall x \varphi(x)$ .

**Lemma 4.5.** *Under these assumptions,  $\Gamma \Rightarrow \forall x \varphi(x)$  is provable if and only if  $\Gamma \Rightarrow \varphi(x)$  is provable.*

*Proof.* Since  $\Gamma$  contains no free variables, the (if) side is proved by directly applying  $(R\forall)$ . The (only if) side follows by the admissibility of the inversion rule of  $(R\forall)$  (see, e.g., Proposition 3.5.4 of [TS00]).  $\square$

**Definition 4.6.** A formula is called *simple* if it is either quantifier free or a universal formula. A set or multiset of formulas is called *simple* if all its formulas are simple.

**Lemma 4.7.** *A formula occurs in a proof tree of  $\Gamma \Rightarrow \varphi(x)$  is either simple or a substitution to  $\varphi[i]$  for some  $i \in [0, k]$ . Moreover, any antecedent formula must be simple, and any succedent formula is either quantifier-free or a substitution to  $\varphi[i]$  for  $i \in [0, k]$ .*

*Proof.* Assume that it is not the case. Find any lowest (i.e. closest to the root) step  $\frac{\Sigma' \Rightarrow \Pi'}{\Sigma \Rightarrow \Pi}$  such that a formula in  $\Sigma', \Pi'$  violates the property above. By checking the eight rules of G3c it's easy to find a contradiction.  $\square$

**Replacing  $(R\rightarrow)$  with  $(R\rightarrow_c)$ .** For technical reasons, we will replace  $(R\rightarrow)$  with the following rule  $(R\rightarrow_c)$  with built-in contradiction, i.e., we introduce an additional Kleene copy  $\alpha \rightarrow \beta$  in the succedent of the premise.

$$\frac{\Gamma, \alpha \Rightarrow \Delta, \beta, \alpha \rightarrow \beta}{\Gamma \Rightarrow \Delta, \alpha \rightarrow \beta} (R\rightarrow_c).$$

**Lemma 4.8.** *Let  $\Gamma, \Delta$  be arbitrary multisets of formulas. Assume that there is a G3c proof tree  $T$  of  $\Gamma \Rightarrow \Delta$ . Then there is a proof tree  $T'$  of  $\Gamma \Rightarrow \Delta$  that has no  $(R\rightarrow)$  and may contain  $(R\rightarrow_c)$ .*

*Proof.* We perform induction on the number of the applications of  $(R\rightarrow)$  in the proof tree. Assume that there is at least one application of  $(R\rightarrow)$ . Let  $(X)$  be an arbitrary application of  $(R\rightarrow)$ . We will replace  $(X)$  to be an application of  $(R\rightarrow_c)$  by adding  $\alpha \rightarrow \beta$  to the succedent of every sequent within the sub-tree rooted at  $(X)$ . It is easy to see that the new proof tree is also a valid proof tree, where  $\alpha \rightarrow \beta$  is either a side formula or a weak formula inside the sub-tree above  $(X)$ . Moreover, the number of  $(R\rightarrow)$  is decreased by 1 so we can complete the proof by the induction hypothesis.  $\square$

<sup>21</sup>Strictly speaking, we define  $\varphi[i]$  as the sub-formula starting from the  $(i + 1)$ -th existential quantifier, where the bounded variables  $y_1, \dots, y_j$  and  $x_1, \dots, x_j$  are replaced by fresh free variables. Here, we slightly abuse the notation to call the fresh variables still by  $y_1, \dots, y_j$  and  $x_1, \dots, x_j$ . The readers should keep in mind that a renaming has been done here to avoid naming conflicts between these bounded variables and other free variables.

**Limiting the applications of  $(R\exists)$  and  $(R\forall)$ .** We also need to perform some proof transformation to “pair” the applications of  $(R\exists)$  and  $(R\forall)$ : we want to have a one-to-one correspondence between the applications of  $(R\exists)$  and  $(R\forall)$  such that the application of  $(R\exists)$  immediately follows from its corresponding application of  $(R\forall)$ . Intuitively, this is possible since the succedent formula of the conclusion is highly structured and the antecedent contains only simple formulas.

**Definition 4.9.** An application  $(X)$  of  $(R\forall)$  is said to be *limited* if it is followed by an application of  $(R\exists)$  such that the principal formula of the  $(R\forall)$  is the same as the active formula of the  $(R\exists)$ . It is said to be *unlimited* if it is not limited.

An application  $(Y)$  of  $(R\exists)$  is said to be *limited* if it immediately follows a limited application of  $(R\forall)$ . It is said to be *unlimited* if it is not limited.

**Definition 4.10.** Let  $T$  be a proof tree,  $x$  be a variable, and  $t$  be a term. The substitution  $T(x/t)$  is defined as the tree where all the free occurrence of  $x$  are replaced by  $t$ .

**Lemma 4.11** ( $\beta$ -substitution). *If  $T$  is a proof tree of  $\Gamma \Rightarrow \Delta$  in  $G3c + (R \rightarrow_c)$ , then  $T(x/t)$  is a proof tree of  $\Gamma(x/t) \Rightarrow \Delta(x/t)$ , provided that  $t$  contains no free occurrence of  $x$  and free variables in  $T, \Gamma, \Delta$ .*

*Proof.* Induction on the proof tree (see, e.g., Lemma 3.5.2 of [TS00]). □

**Lemma 4.12.** *If  $\Gamma \Rightarrow \varphi(x)$  is provable, then there is a proof tree  $T$  of  $\Gamma \Rightarrow \varphi(x)$  such that: (1)  $T$  has no  $(R \rightarrow)$  and may contain  $(R \rightarrow_c)$ ; (2)  $T$  has no unlimited application of  $(R\forall)$ .*

*Proof.* Let  $T$  be a proof tree of  $\Gamma \Rightarrow \varphi(x)$ . By Lemma 4.8, we can assume without loss of generality that  $T$  contains no  $(R \rightarrow)$  and may contain  $(R \rightarrow_c)$ . We will then construct a new proof tree  $T'$  satisfies the property (2) above by removing all the unlimited applications of  $(R\forall)$ . Let

$$\frac{\Gamma \Rightarrow \Delta, \alpha(x/w)}{\Gamma \Rightarrow \Delta, \forall x \alpha(x)} (X)$$

be an uppermost (i.e. farthest to the root) invalid application of  $(R\forall)$ . By Lemma 4.7 we know that  $\alpha$  is a substitution of  $\varphi[i]$  for some  $i \in [1, k]$ , and  $x$  is  $x_i$ . It is easy to see that on the path from  $(X)$  to the root, there must be an application of  $(R\exists)$  with active formula  $\forall x \alpha(x)$ , such that the premise of this  $(R\exists)$  is not derived immediately from an application of  $(R\forall)$  with  $\forall x \alpha(x)$  being principal.<sup>22</sup>

Find an uppermost application of  $(R\exists)$  satisfying the property above. It is of form

$$\frac{\Sigma \Rightarrow \Pi, \forall x_i \alpha(x_i), \exists y_i \forall x_i \alpha'(x_i, y_i)}{\Sigma \Rightarrow \Pi, \exists y_i \forall x_i \alpha'(x_i, y_i)} (Y)$$

Note that  $x_i$  is  $x$  and  $\alpha$  is a substitution of  $\alpha'$  on the existentially bounded variable  $y_i$ . Find all the unlimited rule applications  $(X_1), (X_2), \dots, (X_m)$  of  $(R\forall)$  above  $(Y)$  such that for all  $j \in [m]$ , (1)  $\forall x_i \alpha(x_i)$  is principal in  $(X_j)$ , and (2) there is no rule application of  $(R\forall)$  on the path from  $(X_j)$  to  $(Y)$  that satisfies (1). Clearly  $m \geq 1$  since  $(X)$  satisfies (1).

- We add a limited application of  $(R\forall)$  immediately above  $(Y)$ , as follows.

---

<sup>22</sup>Assume not, this  $\forall x \alpha(x)$  can only be a side formula or an active formula on the path from  $(X)$  to the root. It cannot be an active formula since no rule other than  $(R\exists)$  can has universally quantified active formula on the succedent (recall that we assume all the propositional connectives appear within quantifiers). Therefore it must be side formulas, meaning that it must appear in the succedent of the conclusion, which is impossible.



$$\frac{\frac{\Sigma \Rightarrow \Pi, \alpha(x_i/z), \exists y_i \forall x_i \alpha'(x_i, y_i)}{\Sigma \Rightarrow \Pi, \forall x_i \alpha(x_i), \exists y_i \forall x_i \alpha'(x_i, y_i)} (X')}{\Sigma \Rightarrow \Pi, \exists y_i \forall x_i \alpha'(x_i, y_i)} (Y)$$

where  $z$  is a fresh variable that has no free occurrence in the proof tree. This new formula  $\alpha(x_i/z)$  is added as a side or weak formula into all the sequents in the subtree rooted at  $(Y)$  except for the parts above  $(X_j)$  for some  $j \in [1, m]$ .

- For each  $(X_j)$  with premise  $\Gamma_j \Rightarrow \Delta_j, \phi(x/w_j)$ , we apply Lemma 4.11 to obtain a proof tree of  $\Gamma_j \Rightarrow \Delta_j, \alpha(x_i/z)$  ( $(\alpha(x_i/y_j)(y_j/z) = \alpha(x_i/z))$ ). We then remove the application  $(X_j)$  directly and link this new proof tree back. This is possible since  $\alpha(x_i/z)$  has been available as a side formula.

By doing this we can remove at least one unlimited application of  $(R\forall)$ . The lemma then follows by a simple induction.  $\square$

For convenience, we will use  $(R\exists\forall)$  to denote a pair of adjacent limited  $(R\forall)$  and  $(R\exists)$ , and regard  $(R\exists\forall)$  as a single rule of the following form in a proof tree with only limited applications of  $(R\forall)$  and  $(R\exists)$ .

$$R\exists\forall: \frac{\Gamma \Rightarrow \Delta, \alpha(y/t, x/z), \exists y \forall x \alpha(y, x)}{\Gamma \Rightarrow \Delta, \exists y \forall x \alpha(y, x)}$$

where  $t$  is a term, and  $z$  is a fresh variable without free occurrence in  $\Gamma, \Delta$ , and  $\forall x \alpha(y/t, x)$ . We call  $z$  the *proper variable* in the application of this rule.

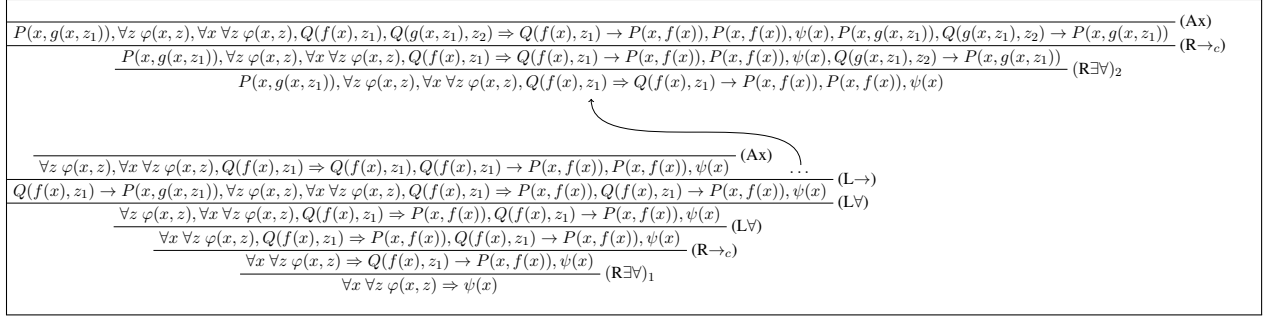
**Canonical Proofs.** Now we show that if  $\Gamma \Rightarrow \varphi(x)$  is provable, then it has a well-structured proof. We need to use the following  $\beta$ -substitution lemma for  $G3c + (R\exists\forall) + (R \rightarrow_c)$ . Similar to previous definition, we define  $T(x/t)$  be the proof tree obtained from  $T$ , where all free occurrences of  $x$  are replaced by  $t$ . It can be easily proved by structural induction on the proof tree  $T$ .

**Lemma 4.13** ( $\beta$ -substitution for  $G3c + (R\exists\forall) + (R \rightarrow_c)$ ). *Assume that  $T$  is a proof tree of  $\Gamma \Rightarrow \Delta$  with  $G3c + (R\exists\forall) + (R \rightarrow_c)$ , then  $T(x/t)$  is a proof tree of  $\Gamma(x/t) \Rightarrow \Delta(x/t)$ , provided that  $t$  contains no free occurrence of  $x$  and free variables in  $T, \Gamma, \Delta$ .*

**Definition 4.14.** A proof tree  $T$  for  $\Gamma \Rightarrow \varphi(x)$  is said to be *canonical* if it satisfies the following three properties:

- (i)  $T$  does not contain  $(R\forall)$ ,  $(R\exists)$ ,  $(R\rightarrow)$ , and  $(L\exists)$ . It may contain  $(R\exists\forall)$  and  $(R\rightarrow_c)$ .
- (ii) Each sequent in  $T$  is of form  $\Gamma \Rightarrow \Phi, \Theta$  where  $\Gamma$  contains only simple formulas,  $\Phi$  contains substitutions to  $\varphi[i]$  for some  $i \in [0, k]$ , and  $\Theta$  contains quantifier-free formulas.
- (iii) Each free variable occurs in  $T$  except for  $x$  is introduced as the proper variable in an application of  $(R\exists\forall)$  for the universal quantifier.

**Theorem 4.15.** *Let  $\varphi(x) = \exists y_1 \forall x_1 \dots \exists y_k \forall x_k \phi(x, \vec{x}, \vec{y})$  be an implicitly bounded formula with exactly one free variable  $x$  for quantifier-free  $\phi$ , and  $\Gamma$  be a set of universal formulas written explicitly in the form  $\forall \vec{x} \beta(\vec{x})$  for some quantifier-free  $\beta(\vec{x})$ . If  $\Gamma \Rightarrow \varphi(x)$  is provable, then there is a canonical proof tree for  $\Gamma \Rightarrow \varphi(x)$ .*



**Figure 2:** A canonical proof tree of  $\forall x \forall z \varphi(x, z) \Rightarrow \psi(x)$ , where  $\varphi(x, z) \triangleq Q(f(x), z) \rightarrow P(x, g(x, z))$ ,  $\psi(x) \triangleq \exists y \forall z (Q(y, z) \rightarrow P(x, y))$ .

*Proof.* We first apply Lemma 4.12 to remove all the unlimited applications of  $(R\forall)$  and replace  $(R\rightarrow)$  by  $(R\rightarrow_c)$ . We add a useless limited  $(R\forall)$  above every unlimited application of  $(R\exists)$  to make it limited (the active formula of this  $(R\forall)$  is added as a side or weak formula in the sequents above). Then we replace all the pairs of adjacent applications of  $(R\forall)$  and  $(R\exists)$  by  $(R\exists\forall)$ .

The second property can be proved similar to Lemma 4.7. The only thing to check is that the new rules  $(R\rightarrow_c)$  and  $(R\exists\forall)$  maintain the subformula property.

For the third property, we only need to consider  $(R\exists\forall)$  and  $(L\forall)$  where a term  $t$  (which may contain new variables) is introduced for the existential and universal quantifiers, respectively. For every such application, if the term  $t$  contains other free variables, we apply Lemma 4.13 to substitute all other free variables to be 0.  $\square$

#### 4.4 Unbounded tree exploration games

Due to technical reasons, we need to define unbounded variants of the tree exploration games and partial game trees. Let  $\varphi(x) = \exists y_1 \forall x_1 \dots \exists y_k \forall x_k \phi(x, \vec{x}, \vec{y})$  be an implicitly bounded  $\mathcal{L}$ -formula as discussed above and  $(\mathcal{M} = (\mathcal{D}, \mathcal{I}), n_0)$  be a board. An *unbounded partial game tree* for  $\varphi$  on the board  $(\mathcal{M}, n_0)$  is defined as a finite rooted tree on which each edge is labeled with a pair  $(m, n) \in \mathcal{D} \times \mathcal{D}$  of elements. The *unbounded tree exploration game* starting from an unbounded partial game tree  $T$  of  $\varphi$  is defined as follows. In each round, the truthifier chooses a node  $u$  on the tree and specifies a number  $m \in \mathcal{D}$ ; the falsifier then specifies a number  $n \in \mathcal{D}$ ; after this round, a child of  $u$  is added to the tree by an edge labeled  $(m, n)$ . The truthifier wins if and only if there is a node on the tree such that the pairs on the path from the root to the node form a satisfying assignment of  $\phi(x/n_0, \vec{x}, \vec{y})$  within  $\mathcal{M}$ , where the truthifier's moves are for  $\vec{y}$  and the falsifier's moves are for  $\vec{x}$ . In particular, the unbounded tree exploration game of  $\varphi$  is defined as the game starting from a tree with only the root.

An  $\mathcal{L}$ -strategy of the truthifier of length  $\ell \in \mathbb{N}$  and initial tree size  $d$  is a sequence

$$\tau = \langle p_1, r_1, p_2, r_2, \dots, p_\ell, r_\ell \rangle,$$

where  $p_i$  is an  $\mathcal{L}$ -term and  $r_i \in \mathbb{N}$  such that  $1 \leq r_i < d + i$ . Let  $(\mathcal{M}, n_0)$  be a board and  $T$  be an unbounded partial game tree on the board with  $d$  nodes numbered  $1, 2, \dots, d$ . The game-theoretic strategy for the unbounded tree exploration game starting from  $T$  induced by  $\tau$  is the following strategy:

- In the  $i$ -th move, the truthifier introduces a node numbered  $d + i$  as a child of the node  $r_i$ , and chooses the element  $v_i \triangleq p_i^{\mathcal{M}}(n_0, T, \Gamma) \in \mathcal{M}$ , where  $\Gamma$  describes the moves of previous rounds (including  $v_1, \dots, v_{i-1}$  and the falsifier's moves).

A length- $\ell$   $\mathcal{L}$ -strategy is said to be a universal winning strategy if the truthifier playing the induced game-theoretic strategy wins within  $\ell$  moves against any strategy of the falsifier on any board  $(\mathcal{M}, n_0)$ .

**Lemma 4.16.** *Let  $\mathcal{T}$  be a bounded theory over the language  $\mathcal{L}$  that is closed under if-then-else. If there is an  $O(1)$ -length  $\mathcal{L}$ -strategy that is a universal winning strategy of the truthifier for the unbounded tree exploration game of  $\varphi = [\psi]_{\text{imp}}$ , then there is an  $O(1)$ -length  $\mathcal{L}$ -strategy that is a universal winning strategy of the truthifier for the tree exploration game of  $\psi$ .*

*Proof.* Assume that  $\tau = \langle p_1, r_1, p_2, r_2, \dots, p_\ell, r_\ell \rangle$  is a universal winning  $\mathcal{L}$ -strategy of length  $\ell \in \mathbb{N}$  for the unbounded tree exploration game. Let  $p'_i(x, \Gamma)$  be the term defined as follows:

- (i) Parse  $\Gamma = (m_1, n_1, m_2, n_2, \dots, m_{i-1}, n_{i-1})$  as the moves in previous rounds.
- (ii) Define  $\Gamma_0$  to be the empty list and  $\hat{\Gamma}_{j+1}$  to be

$$\hat{\Gamma}_{j+1} = \begin{cases} \hat{\Gamma}_j; (m_{j+1}, n_{j+1}) & \text{if } m_{j+1} = p_{j+1}(n_0, \hat{\Gamma}_j) \\ \hat{\Gamma}_j; (p_{j+1}(n_0, \hat{\Gamma}_j), 0) & \text{otherwise} \end{cases}$$

- (iii) Output 0 if  $p_i(n_0, \hat{\Gamma}_{i-1})$  is not a valid move (i.e., it violates the inequality for the bounded variable); and output  $p_i(n_0, \hat{\Gamma}_{i-1})$  otherwise.

Note that such  $p'_i$  always exists as  $\mathcal{T}$  is closed under if-then-else. We now argue that the  $\mathcal{L}$ -quasi-strategy  $\tau' \triangleq \langle p'_1, r_1, p'_2, r_2, \dots, p'_\ell, r_\ell \rangle$  is indeed a universal winning  $\mathcal{L}$ -strategy for the tree exploration game of  $\psi$ . Intuitively,  $\tau'$  is the following  $\mathcal{L}$ -quasi-strategy: it simulates  $\tau$  if it gives a valid move; otherwise, it simply outputs 0 and “forgets” the response of the falsifier, pretending that in this round it simulates  $\tau$  and the falsifier's response were 0.

By the definition of  $p'_i$  it is easy to see that  $\tau'$  is an  $\mathcal{L}$ -strategy for the tree exploration game of  $\psi$ , since it will never output an invalid move. Towards a contradiction we assume that it is not a universal winning strategy. In such case, there exist a board  $(\mathcal{M}, n_0)$  and a strategy  $\tau'_f$  for the falsifier that prevents the truthifier from winning within  $\ell$  rounds on the board against the truthifier playing the induced strategy of  $\tau'$ . We now construct a strategy  $\tau_f$  of the falsifier for the unbounded tree exploration game of  $\varphi$  on the board  $(\mathcal{M}, n_0)$  that prevents  $\tau$  from winning within  $\ell$  rounds and thus leads to a contradiction.

- Assume that the moves of  $\tau'_f$  against  $\tau'$  is  $n'_1, n'_2, \dots, n'_\ell$ . In the  $i$ -th move, if the truthifier's move is an invalid move in the (bounded) tree exploration game (i.e., it violates the inequality for the bounded variable), the falsifier chooses  $n_i \triangleq 0$ ; otherwise the falsifier chooses  $n_i \triangleq n'_i$ .

It is easy to check that against this strategy of the falsifier,  $\tau$  cannot win within  $\ell$  rounds. This is because the transcript of  $\tau_f$  vs  $\tau$  is exactly the lists  $\hat{\Gamma}$  in the definition of  $\tau'$ ; and since  $\tau'$  cannot win against  $\tau'_f$  within  $\ell$  rounds,  $\tau$  also cannot win against  $\tau_f$  within  $\ell$  rounds.  $\square$

This lemma shows that to obtain a winning strategy of the tree exploration game, we only need to extract a winning strategy of the unbounded tree exploration game from a G3c proof tree. In practice, this means that we do not need to treat bounded quantifiers in a special way.

## 4.5 Partial game trees from sequents

Let  $T$  be a canonical proof tree of  $\Gamma \Rightarrow \varphi(x)$ . An  $\mathcal{L}$ -term partial game tree  $T_{\mathcal{L}}$  of  $\varphi(x)$  is a rooted tree that satisfies the following properties. Assume that the nodes are  $v_1, v_2, v_3, \dots, v_\ell$ , where  $v_1$  is the root,  $p_i$  be the parent of  $v_i$ , and  $e_i$  be the edge connecting  $v_i$  and  $p_i$  for  $2 \leq i \leq \ell$ . Let  $\text{path}(i)$  denote the edges on the path from  $v_i$  to the root. Then:

- (i) Each node  $v_i$  is marked as a number  $o_i \in \{1, \dots, \ell\}$  such that  $o_i > o_{p_i}$ , i.e., it is a topological order. This  $o_i$  indicates the order of the introduction of the nodes in the partial game tree.
- (ii) Each edge  $e_i$  is associated with a pair  $(q_i, z_i)$ , where  $z_i$  is a variable and  $q_i$  is an  $\mathcal{L}$ -term. The variables  $\{z_i\}$  are pairwise distinct, and the term  $q_i$  only contains  $x$  and the variables  $z_j$  for  $o_j < o_i$  as free variables.

We also denote it by  $T_{\mathcal{L}}(x, z_1, z_2, \dots, z_\ell)$  to explicitly display the free variables.

Let  $(\mathcal{M} = (\mathcal{D}, \mathcal{I}), n_0)$  be a board. Given an  $\mathcal{L}$ -term partial game tree  $T_{\mathcal{L}}$  and an assignment for the free variables on the tree:  $n_0 \in \mathcal{D}$  for  $x$  and  $n_1, n_2, \dots, n_\ell \in \mathcal{D}$  for the variables  $z_1, z_2, \dots, z_\ell$  on the tree with concrete numbers, we can construct a corresponding *unbounded* partial game tree by replacing  $z_i$  with  $n_i \in \mathcal{D}$  and  $q_i$  with  $q_i^{\mathcal{M}}(x/n_0, z_1/n_1, \dots, z_\ell/n_\ell) \in \mathcal{D}$  for every  $i \in [\ell]$ . We denote this unbounded partial game tree as  $T_{\mathcal{L}}(x/n_0, z_1/n_1, \dots, z_\ell/n_\ell)$ .

We will define a translation  $[\cdot]_{\text{pgt}}$  that maps a sequent in a canonical proof tree  $T$  of  $\Gamma \Rightarrow \varphi(x)$  to a  $\mathcal{L}$ -term partial game tree.<sup>23</sup> To gain some intuition, it is instructive to image the (canonical) proof tree as formed in a proof search procedure, that is, from the root to the leaves. The substitutions of  $\varphi[i]$  in the succedent are introduced as the active formulas when we apply  $(R\forall\exists)$ . At the same time, the principal formula in the conclusion occurs as the Kleene copy in the premise. For every sequent  $\Gamma \Rightarrow \Phi, \Theta$  as described above, the formulas in  $\Phi$  were introduced in the proof search as a tree structure, where the initial formula  $\varphi(x)$  at the root of the proof tree is the root, and the active formula introduced when applying  $(R\forall\exists)$  is a child of the principal formula (which becomes the Kleene copy in the premise) of this rule application. This tree structure naturally corresponds to an  $\mathcal{L}$ -term partial game tree  $T_G = [\Gamma \Rightarrow \Phi, \Theta]_{\text{pgt}}$ , where the proper variables occur as  $z_j$  in  $T_G$  and the witnessing terms occur as  $q_j$  in  $T_G$ .

Formally, we defined the set of formulas *obtained from substitutions of  $\varphi[\cdot]$* , denoted by  $\text{Sub}(\varphi)$ , as follows.<sup>24</sup> A formula  $\psi$  in the succedent of the canonical proof tree is in  $\text{Sub}(\varphi)$  if one of the following conditions hold:

- (i) The sequent containing  $\psi$  is the conclusion and  $\psi$  is  $\varphi(x)$ ;
- (ii) In the rule application with the sequent containing  $\psi$  being the premise,  $\psi$  is a side formula, and its counterpart in the conclusion is in  $\text{Sub}(\varphi)$ ;
- (iii) The sequent containing  $\psi$  is the premise of a rule application of  $(R\rightarrow_c)$ ,  $\psi$  is the Kleene copy, and the principal formula is in  $\text{Sub}(\varphi)$ ;
- (iv) The sequent containing  $\psi$  is the premise of a rule application of  $(R\exists\forall)$ , and  $\psi$  is one of the two active formulas.

<sup>23</sup>Strictly speaking, the translation  $[\cdot]_{\text{pgt}}$  is defined over the occurrences of the sequents in the proof tree  $T$ . In particular, two occurrences of the same sequent may be translated to different partial game trees. Here (and below), we slightly abuse the notation to say  $[\Gamma \Rightarrow \Psi]_{\text{pgt}}$  for a sequent  $\Gamma \Rightarrow \Psi$ , when the occurrence of  $\Gamma \Rightarrow \Psi$  that we are referring to is clear in the context.

<sup>24</sup>Strictly speaking,  $\text{Sub}(\varphi)$  should be a set of occurrences of formulas in the sequent instead of a set of formulas. Here (and below), we slightly abuse the notation to say that a formula is in  $\text{Sub}(\varphi)$ , which actually means that the occurrence of the formula is in  $\text{Sub}(\varphi)$ , when the occurrence is clear in the context.

We can then write each sequent in the proof tree in the form  $\Sigma \Rightarrow \Pi, \Delta$ , where  $\Pi$  contains all the formulas in  $\text{Sub}(\varphi)$ . It then follows from Lemma 4.7 that  $\Delta$  contains only quantifier-free formulas.<sup>25</sup>

Now we define the translation by induction on the depth (i.e., the distance to the root) of the sequent  $\Sigma \Rightarrow \Pi, \Delta$  we want to translate within  $T$ . We will maintain the following properties during this recursive definition.

- **(Structure of Formulas in  $\Pi$ .)** Let  $\Sigma \Rightarrow \Pi, \Delta$  be a sequent in  $T$  where  $\Pi$  contains formulas in  $\text{Sub}(\varphi)$  and  $\Delta$  contains quantifier-free formulas. For all  $\alpha \in \Pi$ , there is an  $i \in [0, k]$  such that

$$\alpha = \exists y_{i+1} \forall x_{i+1} \dots \exists y_k \forall x_k \phi(x, x_1/z_1, \dots, x_i/z_i, y_1/q_1, \dots, y_i/q_i),$$

where  $q_j$  is an  $\mathcal{L}$ -term with no free variables except for  $x$  and the free variables on the  $\mathcal{L}$ -term partial game tree corresponding to the parent of this sequent (if it has a parent).

- **(Formulas in  $\Pi \sim$  Nodes in Partial Game Tree).** For each sequent  $\Sigma \Rightarrow \Pi, \Delta$  in  $T$  as described above, there is a one-to-one correspondence between the formulas  $\alpha \in \Pi$  and the nodes  $v \in [\Sigma \Rightarrow \Pi, \Delta]_{\text{pgt}}$ . Let  $\alpha \in \Pi$  be a formula of the form

$$\alpha = \exists y_{i+1} \forall x_{i+1} \dots \exists y_k \forall x_k \phi(x, x_1/z_{j_1}, \dots, x_i/z_{j_i}, y_1/q_{j_1}, \dots, y_i/q_{j_i}).$$

Then the edges on the path from the root to the node  $v$  corresponding to  $\alpha$  is associated with the pairs

$$(q_{j_1}, z_{j_1}), (q_{j_2}, z_{j_2}), \dots, (q_{j_{i-1}}, z_{j_{i-1}}).$$

Let  $\Sigma \Rightarrow \Pi, \Delta$  be a sequent in  $T$  such that  $\Pi$  contains formulas in  $\text{Sub}(\varphi)$  and  $\Delta$  contains quantifier-free formulas.

- If the depth is 0, (i.e., we are at the root of the proof tree), we know that  $\Delta = \emptyset$  and  $\Pi = \{\varphi(x)\}$ . We define  $[\Sigma \Rightarrow \Pi, \Delta]_{\text{pgt}}$  as a  $\mathcal{L}$ -term partial game tree with only the root.
- Otherwise we assume that the parent of the sequent is  $\Sigma' \Rightarrow \Pi', \Delta'$ , where  $\Pi'$  contains formulas in  $\text{Sub}(\varphi)$  and  $\Delta'$  contains quantifier-free formulas. We consider which rule is applied to derive  $\Sigma' \Rightarrow \Pi', \Delta'$  from  $\Sigma \Rightarrow \Pi, \Delta$  (and another premise, if  $(L \rightarrow)$  is applied). There are the following possibilities:  $(L \rightarrow)$ ,  $(R \rightarrow_c)$ ,  $(L \forall)$ , and  $(R \exists \forall)$ .

(i) If the rule applied is one of  $(L \rightarrow)$ ,  $(R \rightarrow_c)$ , and  $(L \forall)$ , we know that  $\Pi = \Pi'$ , and we define  $[\Sigma \Rightarrow \Pi, \Delta]_{\text{pgt}} \triangleq [\Sigma' \Rightarrow \Pi', \Delta']_{\text{pgt}}$ .

(ii) If the rule applied is  $(R \exists \forall)$  with  $\exists y_i \forall x_i \alpha(x_i, y_i)$  being principal, we know that  $\Pi = \Pi' \cup \{\alpha(x_i/z, y_i/t)\}$  for some fresh variable  $z$  and term  $t$ . By the definition of canonical proof tree and the invariant, we know that  $t$  has no free variables except for  $x$  and free variables in  $\Sigma', \Pi', \Delta'$ , which further means that  $t$  has no free variables except for  $x$  and the free variables in  $[\Sigma' \Rightarrow \Pi', \Delta']_{\text{pgt}}$ .

Let  $T'_{\text{pgt}} \triangleq [\Sigma' \Rightarrow \Pi', \Delta']_{\text{pgt}}$  be a partial game tree with  $\ell$  nodes (including the root). Let  $v \in T'_{\text{pgt}}$  be the node corresponding to the formula  $\exists y_i \forall x_i \alpha(x_i, y_i) \in \Pi'$ , as promised in the induction

<sup>25</sup>Recall that the conclusion of the proof tree is  $\Gamma \Rightarrow \varphi(x)$ , where  $\Gamma$  contains universal sentences written explicitly in the form  $\forall \vec{z} \beta$  for some quantifier-free  $\beta$ . A formula  $\alpha$  in the proof tree that is not in  $\text{Sub}(\varphi)$  is necessarily a sub-formula of some  $\forall \vec{z} \beta \in \Gamma$ . As it appears in the succedent, on the path from it to the root there must be an application of  $(L \rightarrow)$  that “carries” it from the antecedent to the succedent. Since all the propositional connectives appears within the quantifier-free  $\beta$ , this formula  $\alpha$  has to be quantifier-free.

hypothesis. We add a child of  $v$  labeled  $\ell + 1$  that is connected to  $v$  with an edge associated with the pair  $(t, z)$ . This child corresponds to the active formula  $\alpha(x_i/z, y_i/t) \in \Gamma \setminus \Gamma'$ , whereas the correspondence between other nodes and formulas are induced from the correspondence in  $T'_{\text{pgt}}$ . (Note that the formulas in  $\Pi'$  appears as side formulas or the Kleene copy in  $\Pi$ .) It's easy to verify that the two properties for the induction is satisfied and the new tree is still a valid  $\mathcal{L}$ -term partial game tree.

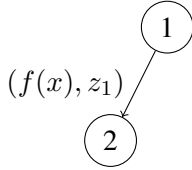
*Example 4.17.* For instance, consider the canonical proof tree of the sequent  $\forall x \forall z (Q(f(x), z) \rightarrow P(x, g(x, z))) \Rightarrow \exists y \forall z (Q(y, z) \rightarrow P(x, y))$  in Figure 2. The  $\mathcal{L}$ -term partial game tree corresponding to the sequent on the root is a tree with only the root; the game trees corresponding to the premises of  $(R\exists\forall)_1$  and  $(R\exists\forall)_2$  are shown in Figure 3a and Figure 3b, respectively. Node 1, Node 2, and Node 3 correspond to the following substitutions of  $\exists y \forall z (Q(y, z) \rightarrow P(x, y))$  in the proof tree:

Node 1:  $\exists y \forall z (Q(y, z) \rightarrow P(x, y))$ .

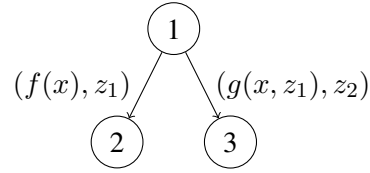
Node 2:  $Q(f(x), z_1) \rightarrow P(x, f(x))$ .

Node 3:  $Q(g(x, z_1), z_2) \rightarrow P(x, g(x, z_1))$ .

The partial game trees corresponding to the other sequents in the proof tree are inherited from the partial game trees of their parents.



(a) PGT of the premise of  $(R\exists\forall)_1$ .



(b) PGT of the premise of  $(R\exists\forall)_2$ .

**Figure 3:** Partial game trees corresponding to sequents in the proof tree in Figure 2.

## 4.6 Completing the argument: Winning strategies from proofs

**Definition 4.18.** Let  $\tau = \langle p_1, r_1, p_2, r_2, \dots, p_\ell, r_\ell \rangle$  and  $\tau' = \langle p'_1, r'_1, p'_2, r'_2, \dots, p'_{\ell'}, r'_{\ell'} \rangle$  be two  $\mathcal{L}$ -term strategies with initial tree size  $d$ . We define the concatenation of the two strategies as

$$\tau; \tau' \triangleq \langle p_1, r_1, p_2, r_2, \dots, p_\ell, r_\ell, p'_1, r'_1, p'_2, r'_2, \dots, p'_{\ell'}, r'_{\ell'} \rangle.$$

Note that  $\tau; \tau'$  is an  $\mathcal{L}$ -term strategy with length  $\ell + \ell'$  and initial tree size  $d$ .

**Lemma 4.19.** Let  $T$  be a canonical proof tree of  $\Gamma \Rightarrow \varphi(x)$ ,  $\Sigma \Rightarrow \Pi$ ,  $\Delta$  be a sequent in  $T$  (where  $\Pi$  contains formulas in  $\text{Sub}(\varphi)$  and  $\Delta$  contains quantifier-free formulas), and  $T_G(x, z_1, z_2, \dots, z_d) \triangleq [\Sigma \Rightarrow \Pi, \Delta]_{\text{pgt}}$ . There is an  $O(1)$  length  $\mathcal{L}$ -term strategy  $\tau$  with initial tree size  $d + 1$  such that the following holds:

- **(Correctness).** Let  $(\mathcal{M}, n_0)$  be a board and  $\sigma$  be an assignment of free variables in this  $\Sigma \Rightarrow \Pi, \Delta$  such that  $\sigma(x) = n_0$ . If in the model  $\mathcal{M}$  and with the assignment  $\sigma$  to free variables, every formula in  $\Sigma$  is true and every formula in  $\Delta$  is false, the induced game-theoretic strategy of  $\tau$  is a winning strategy for the unbounded tree exploration game starting from  $T_G(x/n_0, z_1/\sigma(z_1), z_2/\sigma(z_2), \dots, z_d/\sigma(z_d))$ .
- **(Extraction).** There is a polynomial-time algorithm that outputs the strategy for every sequent in  $T$  given the canonical proof tree  $T$  as input.



*Proof.* Let  $\Sigma \Rightarrow \Pi, \Delta$  be a sequent in  $\mathcal{T}$ , we prove that if this lemma holds for the premise(s) of this  $\Sigma \Rightarrow \Pi, \Delta$ , then it also holds for this  $\Sigma \Rightarrow \Pi, \Delta$ . The lemma then follows from a structural induction on  $T$ . Consider the rule application on this  $\Sigma \Rightarrow \Pi, \Delta$ .

**Case 1.** If it is an application of (Ax), we consider the following two cases.

**Case 1.1.** Assume that the principal formula  $P \notin \text{Sub}(\varphi)$ , i.e., the sequent is of form  $\Sigma', P \Rightarrow \Pi, \Delta', P$  for some atomic predicate  $P$ , where  $\Sigma = \Sigma', P$  and  $\Delta = \Delta', P$ . Then for every board  $(\mathcal{M}, n_0)$  and assignment  $\sigma$  for free variables in  $\Sigma \Rightarrow \Pi, \Delta$ , in the model  $\mathcal{M}$  and with the assignment  $\sigma$ , it is impossible for all the formulas in  $\Sigma$  to be true and all the formulas in  $\Pi$  to be false, since  $P \in \Sigma \cap \Pi$ .

**Case 1.2.** Otherwise, the principal formula  $P \in \text{Sub}(\varphi)$ , i.e., the sequent is of form  $\Sigma', P \Rightarrow \Pi', \Delta, P$ , where  $\Sigma = \Sigma', P$  and  $\Delta = \Delta', P$ . Let  $T_G \triangleq [\Sigma \Rightarrow \Pi, \Delta]_{\text{pgt}}$ . Let  $(\mathcal{M}, n_0)$  be a board and  $\sigma$  be an assignment of free variables in  $\Sigma \Rightarrow \Pi, \Delta$  such that  $\sigma(x) = n_0$ . In the model  $\mathcal{M}$  and with the assignment  $\sigma$  to free variables, if all the formulas in  $\Sigma$  are true, then  $P \in \Sigma$  is true, meaning that the node in  $T_G$  corresponding to  $P \in \Pi$  is a winning node of the evaluation game of  $\varphi$ . In such case, a trivial strategy that does nothing is a winning strategy.

**Case 2.** If it is an application of (L $\perp$ ), the sequent is of form  $\Sigma', \perp \Rightarrow \Pi, \Delta$ . Similar to Case 1, it is impossible for all the formulas in  $\Sigma$  to be true, since  $\perp \in \Sigma$ .

**Case 3.** If it is an application of (L $\rightarrow$ ), it is of form

$$\frac{\Sigma' \Rightarrow \Pi, \Delta, \alpha \quad \Sigma', \beta \Rightarrow \Pi, \Delta}{\Sigma', \alpha \rightarrow \beta \Rightarrow \Pi, \Delta}$$

where  $\Sigma = \Sigma', \alpha \rightarrow \beta$ . Let  $T_G \triangleq [\Sigma' \Rightarrow \Pi, \Delta, \alpha]_{\text{pgt}} = [\Sigma', \beta \Rightarrow \Pi, \Delta]_{\text{pgt}} = [\Sigma \Rightarrow \Pi, \Delta]_{\text{pgt}}$ . The formula  $\alpha \rightarrow \beta$  is simple. By the induction hypothesis, we know that there are  $\mathcal{L}$ -strategies  $\tau_1$  for  $\Sigma' \Rightarrow \Pi, \Delta, \alpha$  and  $\tau_2$  for  $\Sigma', \beta \Rightarrow \Pi, \Delta$  satisfying this lemma. We know prove that  $\tau_1; \tau_2$  is the desired strategy for this  $\Sigma \Rightarrow \Pi, \Delta$ . Let  $(\mathcal{M}, n_0)$  be a board and  $\sigma$  be an assignment of free variables such that  $\sigma(x) = n_0$ . If in the model  $\mathcal{M}$  and with the assignment  $\sigma$  to free variables, all formulas in  $\Sigma = \Sigma', \alpha \rightarrow \beta$  are true and all formulas in  $\Delta$  are false, since  $(\alpha \rightarrow \beta) \in \Sigma$  either  $\alpha$  is false or  $\beta$  is true. We can see from the induction hypothesis that: In the former case,  $\tau_1$  will be a winning strategy for the unbounded tree exploration game starting from  $T_G(x/n_0, z_1/\sigma(z_1), \dots, z_\ell/\sigma(z_\ell))$ ; in the latter case,  $\tau_2$  will be a winning strategy for the same game. As a result,  $\tau_1; \tau_2$  is always a winning strategy for the unbounded tree exploration game starting from  $T_G(x/n_0, z_1/\sigma(z_1), \dots, z_\ell/\sigma(z_\ell))$ .

**Case 4.** If it is an application of (R $\rightarrow_c$ ), consider the following two cases.

**Case 4.1.** Assume that the principal formula is not in  $\text{Sub}(\varphi)$ . Then the rule applications is of form

$$\frac{\Sigma, \alpha \Rightarrow \Pi, \Delta', \beta, \alpha \rightarrow \beta}{\Sigma \Rightarrow \Pi, \Delta', \alpha \rightarrow \beta}$$

where  $\Delta = \Delta', \beta$ . By the definition, we know that  $\beta$  in the premise is not in  $\text{Sub}(\varphi)$ . Let  $T_G \triangleq [\Sigma, \alpha \Rightarrow \Pi, \Delta', \beta]_{\text{pgt}} = [\Sigma \Rightarrow \Pi, \Delta]_{\text{pgt}}$ . Similar to Case 3, we can see the  $\mathcal{L}$ -term strategy for the premise from the induction hypothesis is indeed the required  $\mathcal{L}$ -term strategy for  $\Sigma \Rightarrow \Pi, \Delta$ , because in any model and for any assignment to free variables, if  $\alpha \rightarrow \beta$  is false, then  $\alpha$  is true and  $\beta$  is false.

**Case 4.2.** Now we assume that the principal formula is in  $\text{Sub}(\varphi)$ . It is of form

$$\frac{\Sigma, \alpha \Rightarrow \Pi', \Delta, \beta, \alpha \rightarrow \beta}{\Sigma \Rightarrow \Pi', \Delta, \alpha \rightarrow \beta}$$

where  $\Pi = \Pi', \alpha \rightarrow \beta$ . Let  $T_G \triangleq [\Sigma \Rightarrow \Pi, \Delta]_{\text{pgt}} = [\Sigma, \alpha \Rightarrow \Pi', \Delta, \beta, \alpha \rightarrow \beta]_{\text{pgt}}$ . We claim that the  $\mathcal{L}$ -term strategy for the premise from the induction hypothesis is also the required  $\mathcal{L}$ -term strategy for  $\Sigma \Rightarrow \Pi, \Delta$ . Let  $(\mathcal{M}, n_0)$  be a board and  $\sigma$  be an assignment of free variables in  $\Sigma \Rightarrow \Pi, \Delta$  such that  $\sigma(x) = n_0$ . Assume that in the model  $\mathcal{M}$  and with the assignment  $\sigma$  to free variables, all the formulas in  $\Sigma$  are true and all the formulas in  $\Pi$  are false. Consider the following two cases:

- (i) If  $\alpha \rightarrow \beta$  is true, then the node corresponding to  $\alpha \rightarrow \beta$  in  $T_G$  is a winning node for the truthifier, which means that the truthifier wins the game regardless of its further moves.
- (ii) Otherwise,  $\alpha$  is true and  $\beta$  is false. In such case, we can see from the induction hypothesis that the strategy from the premise will win in the tree exploration game starting from  $T_G(x/n_0, z_1/\sigma(z_1), \dots, z_\ell/\sigma(z_\ell))$  on the board  $(\mathcal{M}, n_0)$ .

As a result, the strategy from the premise is always a winning strategy.

**Case 5.** If it is an application of  $(L\forall)$ , it is of form

$$\frac{\Sigma', \alpha(x/t) \Rightarrow \Pi, \Delta}{\Sigma', \forall x \alpha(x) \Rightarrow \Pi, \Delta}$$

where  $\Sigma = \Sigma', \forall x \alpha(x)$ . Similar to Case 3 and Case 4, we can see that the  $\mathcal{L}$ -term strategy for the premise from the induction hypothesis is indeed the required  $\mathcal{L}$ -term for  $\Sigma \Rightarrow \Pi, \Delta$ , since in any model and for any assignment to free variables, if  $\forall x \alpha(x)$  is true, then  $\alpha(x/t)$  is true.

**Case 6.** If it is an application of  $(R\exists\forall)$ , it is of form

$$\frac{\Sigma \Rightarrow \Pi', \Delta, \alpha(x/z, y/t), \exists y \forall x \alpha(y, x)}{\Sigma \Rightarrow \Pi', \Delta, \exists y \forall x \alpha(y, x)}$$

where  $\Pi = \Pi', \exists y \forall x \alpha(y, x)$  and  $\exists y \forall x \alpha(y, x) \in \text{Sub}(\varphi)$ . Let

$$\begin{aligned} T_G(z_1, z_2, \dots, z_{d-1}) &\triangleq [\Sigma \Rightarrow \Pi, \Delta]_{\text{pgt}}, \\ T'_G(z_1, z_2, \dots, z_{d-1}, z) &\triangleq [\Sigma \Rightarrow \Pi', \Delta, \alpha(x/z, y/t), \exists y \forall x \alpha(y, x)]_{\text{pgt}}, \end{aligned}$$

where  $d$  is the number of nodes in  $T_G$ . Let  $v_i$  be the node corresponding to  $\alpha(x/z, y/t)$  in  $T'_G$ . Assume that  $\exists y \forall x \alpha(y, x)$  corresponds to  $v_j$  in  $T_G$ . By the induction hypothesis, there is a strategy  $\tau_1 = \langle p_1, r_1, p_2, r_2, \dots, p_\ell, r_\ell \rangle$  for the premise that satisfies the lemma. We now show that the strategy  $\tau = \langle t, j, p_1, r_1, p_2, r_2, \dots, p_\ell, r_\ell \rangle$ , which means (intuitively) that the truthifier will choose the node  $v_j$  and add a number  $t$  in the first round and then follow  $\tau_1$ , is the required  $\mathcal{L}$ -term strategy for  $\Sigma \Rightarrow \Pi, \Delta$ .

Let  $(\mathcal{M}, n_0)$  be a board and  $\sigma$  be an assignment of free variables in  $\Sigma \Rightarrow \Pi, \Delta$  such that  $\sigma(x) = n_0$ . Assume that in the model  $\mathcal{M}$  with the assignment  $\sigma$  to free variables, all the formulas in  $\Sigma$  are true and all the formulas in  $\Delta$  are false. We need to prove that the induced game-theoretic strategy of  $\tau$  is a winning strategy of the unbounded tree exploration game starting from  $T_G(x/n_0, z_1/\sigma(z_1), \dots, z_\ell/\sigma(z_\ell))$ .

According to the definition of  $\tau$ , the truthifier will choose  $v_j \in T_G$  (i.e., the node corresponding to  $\exists y \forall x \alpha(y, x)$ ) and a number  $t^M$ . Assume that the falsifier chooses  $n_1 \in \mathcal{M}$  in its move. Let  $\sigma_1 \triangleq \sigma \cup \{z/n_1\}$  be an assignment of free variables in the premise.

In the model  $\mathcal{M}$  and with the assignment  $\sigma_1$ , for the premise of this application of  $(R\exists\forall)$ , all the formulas in the antecedent are true and all the formulas that are not in  $\text{Sub}(\varphi)$  in the succedent are false. Therefore by the induction hypothesis, the induced game-theoretic strategy  $\tau_1$  is a winning strategy for the unbounded tree exploration game starting from  $T'_G$  in the model  $\mathcal{M}$  and with the assignment  $\sigma_1$ . Since  $\tau$  starts to simulate  $\tau_1$  from the second step, the induced game-theoretic strategy of  $\tau$  will also win.

Note that although we do not explicitly discuss the extraction property, it is easy to verify that we can extract the strategy inductively using the proof above.  $\square$

**Theorem** (Reminder of Theorem 4.1). *Let  $\mathcal{T}$  be a universal bounded theory with vocabulary  $\mathcal{L}$  that is closed under if-then-else. Let  $\varphi$  be a bounded  $\mathcal{L}$ -formula of the form*

$$\begin{aligned} \varphi(x) \triangleq & \exists y_1 \leq t_1(x) \forall x_1 \leq s_1(x, y_1) \exists y_2 \leq t_2(x, y_1, x_1) \dots \forall x_{k-1} \leq s_{k-1}(x, y_1, x_1, \dots, y_{k-1}) \\ & \exists y_k \leq t_k(x, y_1, x_1, \dots, y_{k-1}, x_{k-1}) \forall x_k \leq s_k(x, y_1, x_1, \dots, y_k) \phi(x, x_1, \dots, x_k, y_1, \dots, y_k), \end{aligned}$$

where  $\phi(x, \vec{x}, \vec{y})$  is a quantifier-free  $\mathcal{L}$ -formula. Then  $\mathcal{T} \vdash \forall x \varphi(x)$  if and only if there is a universal winning  $\mathcal{L}$ -strategy of length  $O(1)$  for the truthifier in the corresponding tree exploration game of  $\varphi(x)$ .

*Proof.* (If). Assume that  $\mathcal{T} \not\vdash \forall x \varphi(x)$ . Then by the completeness theorem there is a model  $\mathcal{M} = (\mathcal{D}, \mathcal{I})$  and  $n_0 \in \mathcal{D}$  such that  $\varphi^M(n_0)$  is false, which further means that there is a winning strategy of the falsifier in the evaluation game of  $\varphi(x)$  on the board  $(\mathcal{M}, n_0)$ . Consider the strategy of the falsifier in the tree exploration game that simply simulates this winning strategy, i.e., after the truthifier adds a node and specifies an element on the edge, the falsifier treats the path from the root to this node as a partial transcript of the evaluation game and chooses an element according to the strategy of the evaluation game. It is clear that the truthifier cannot reach a winning node, thus it does not have a universal winning  $\mathcal{L}$ -strategy of the tree exploration game.

(Only If). Let  $\mathcal{L}, \mathcal{T}$  and  $\varphi$  be defined as above. Suppose that  $\mathcal{T} \vdash \forall x \varphi(x)$ . By Lemma 4.5 and Theorem 4.15, there is a canonical proof tree for  $\mathcal{T} \vdash \varphi(x)$ . By applying Lemma 4.19 and considering the root of the proof tree, there is an  $O(1)$ -length  $\mathcal{L}$ -strategy  $\tau$  for the unbounded tree exploration game of  $\varphi(x)$ . Finally, we can obtain an  $O(1)$ -length universal winning  $\mathcal{L}$ -strategy  $\tau$  for the tree exploration game of  $\varphi(x)$  by Lemma 4.16.  $\square$

## 4.7 A special case: Falsifiers with oblivious strategies

In this section, we present a special case of game-theoretic witnessing (Theorem 4.1) that involves *sequential* invocations of the *evaluation game* played against an *oblivious falsifier*. This version is sufficient to show the unprovability of strong circuit lower bounds in bounded arithmetic (Section 6.2).

We assume familiarity with the notation from Section 4.1. In particular, let  $\mathcal{T}, \mathcal{L}$ , and  $\varphi(x)$  be defined as in Section 4.1. The main difference is that here we consider the evaluation game (as opposed to the tree exploration game) in the presence of *ancillary information for the truthifier*, as explained next.

**Strategies with Ancillary Information.** Let  $\mathcal{M} = (\mathcal{D}, \mathcal{I})$  be a model for  $\mathcal{T}$ . An  $\mathcal{L}$ -strategy for the truthifier with *ancillary information* in the evaluation game of  $\varphi(x)$  is a sequence  $\tau^t = (p_1, p_2, \dots, p_k)$  of  $k$   $\mathcal{L}$ -terms, where  $p_i \triangleq p_i(n_0, m_1, n_1, \dots, m_{i-1}, n_{i-1}, \vec{a})$  means given the *ancillary information*  $\vec{a}$  (constantly

many elements from  $\mathcal{D}$ ),  $n_0 \in \mathcal{D}$ , and moves  $m_1, n_1, \dots, m_{i-1}, n_{i-1} \in \mathcal{D}$ , the truthifier chooses  $m_i = p_i^{\mathcal{M}}(n_0, m_1, n_1, \dots, m_{i-1}, n_{i-1}, \vec{a})$  as the current move. For every  $\vec{a} \in \vec{\mathcal{D}}$ , the strategy induced by  $\tau^{\mathfrak{t}}$  given  $\vec{a}$  as ancillary information is denoted by  $\tau^{\mathfrak{t}}[\vec{a}]$ . In particular, if the  $\mathcal{L}$ -strategy has no ancillary information, the induced strategy is denoted by  $\tau^{\mathfrak{t}}[\emptyset]$ . Similarly to Section 4.1, the *transcript* of a game given strategies  $\tau^{\mathfrak{t}}$  for the truthifier (possibly with ancillary information) and  $\tau^{\mathfrak{f}}$  for the falsifier, denoted by  $\langle \tau^{\mathfrak{t}} : \tau^{\mathfrak{f}} \rangle$ , is a pair  $(\vec{n}, \vec{m})$  of sequences that records the moves of both players.

**Theorem 4.20** (Winning strategies against oblivious falsifiers). *Let  $\mathcal{T}$  be a universal theory over the language  $\mathcal{L}$  that is closed under if-then-else. Let  $\varphi(x)$  be the formula*

$$\begin{aligned} \varphi(x) \triangleq & \exists y_1 \leq t_1(x) \forall x_1 \leq s_1(x, y_1) \exists y_2 \leq t_2(x, y_1, x_1) \dots \forall x_{k-1} \leq s_{k-1}(x, y_1, x_1, \dots, y_{k-1}) \\ & \exists y_k \leq t_k(x, y_1, x_1, \dots, y_{k-1}, x_{k-1}) \forall x_k \leq s_k(x, y_1, x_1, \dots, y_k) \phi(x, x_1, \dots, x_k, y_1, \dots, y_k), \end{aligned}$$

where  $\phi(x, \vec{x}, \vec{y})$  is a quantifier-free  $\mathcal{L}$ -formula. If  $\mathcal{T} \vdash \forall x \varphi(x)$ , then there is a constant  $\ell \in \mathbb{N}$  and  $\mathcal{L}$ -strategies  $\tau_1^{\mathfrak{t}}, \tau_2^{\mathfrak{t}}, \dots, \tau_\ell^{\mathfrak{t}}$  (with ancillary information) such that, for any board  $(\mathcal{M}, n_0)$  and evaluation game of  $\varphi(x)$  on  $(\mathcal{M}, n_0)$ , for every strategy  $\tau^{\mathfrak{f}}$  of the falsifier:

- either  $\hat{\tau}_1^{\mathfrak{t}} \triangleq \tau_1^{\mathfrak{t}}[\emptyset]$  beats  $\tau^{\mathfrak{f}}$ ,
- or  $\hat{\tau}_2^{\mathfrak{t}} \triangleq \tau_2^{\mathfrak{t}}[\langle \hat{\tau}_1^{\mathfrak{t}} : \tau^{\mathfrak{f}} \rangle]$  beats  $\tau^{\mathfrak{f}}$ ,
- or  $\hat{\tau}_3^{\mathfrak{t}} \triangleq \tau_3^{\mathfrak{t}}[\langle \hat{\tau}_1^{\mathfrak{t}} : \tau^{\mathfrak{f}} \rangle, \langle \hat{\tau}_2^{\mathfrak{t}} : \tau^{\mathfrak{f}} \rangle]$  beats  $\tau^{\mathfrak{f}}$ ,
- ...,
- or  $\hat{\tau}_\ell^{\mathfrak{t}} \triangleq \tau_\ell^{\mathfrak{t}}[\langle \hat{\tau}_1^{\mathfrak{t}} : \tau^{\mathfrak{f}} \rangle, \langle \hat{\tau}_2^{\mathfrak{t}} : \tau^{\mathfrak{f}} \rangle, \dots, \langle \hat{\tau}_{\ell-1}^{\mathfrak{t}} : \tau^{\mathfrak{f}} \rangle]$  beats  $\tau^{\mathfrak{f}}$ .

Before we establish this result, a few comments are in order. First, notice that the moves of the falsifier only depend on the previous moves in the *current game*. On the other hand, the truthifier gets as ancillary information the transcripts of all previous games, and succeeds in beating the strategy of the falsifier after (sequentially) playing at most  $\ell = O(1)$  games. Intuitively, the falsifier is *oblivious*, since its moves in the current game do not depend on the moves from any previously completed or different game played in parallel, as in the tree exploration game described in Section 4.1. Consequently, when extracting computational information from proofs (where one defines appropriate strategies for the falsifier and considers the behaviour of the truthifier), Theorem 4.20 is more limited than Theorem 4.1.

*Proof of Theorem 4.20.* Intuitively, as explained above, the meaning of the theorem is that the truthifier has a winning strategy (with ancillary information) in  $\ell$  sequential plays of the evaluation game when the falsifier's strategy is fixed. We will obtain such a strategy from a strategy for the truthifier that succeeds in the tree exploration game. This is not entirely obvious, since there is a mismatch between the games: the next play of the truthifier in the tree exploration game depends on all previous plays in the game tree, while in the evaluation game there is no game tree and they play a sequence of evaluation games.

Let  $\mathcal{T}, \mathcal{L}$ , and  $\varphi(x)$  be defined as above. By Theorem 4.1, there exists an  $\ell = O(1)$  length  $\mathcal{L}$ -term winning strategy  $\tau^{\text{tree}}$  for the tree exploration game of  $\varphi(x)$ . Let  $(\mathcal{M}, n_0)$  be a board. Consider the tree exploration game when the falsifier plays with a *fixed strategy of the evaluation game*, that is, there exist functions  $f_1(x, y_1), f_2(x, y_1, y_2), \dots, f_k(x, y_1, y_2, \dots, y_k)$  such that:

- In the  $i$ -th step, if the truthifier adds a node  $v$  to the partial game tree as a child of  $u$  and chooses  $m$  as the label and  $(m_1, n_1), (m_2, n_2), \dots, (m_d, n_d)$  are the labels on the length- $d$  path from the root to  $u$ , then the falsifier's move is  $f_{d+1}(n_0, m_1, m_2, \dots, m_d, m)$ .

We say that a falsifier’s strategy of this form in the tree exploration game is *oblivious*, i.e., the next move of the falsifier only considers moves in the corresponding root-to-node path. Since  $\tau^{\text{tree}}$  is a winning strategy for the tree exploration game, it beats all strategies of the falsifier, including oblivious strategies.

We would like to simulate  $\tau^{\text{tree}}$ , a strategy for the tree exploration game, in the context of Theorem 4.20, where the evaluation game is played sequentially and the truthifier has ancillary information. The main idea is to play each *round* in the *tree exploration game* as a new *game* in the *evaluation game* that simulates the current root-to-node path. This guarantees when translating strategies that all the necessary information from the tree exploration game appears in the transcript of previous plays (ancillary information) during the next evaluation game. (If the root-to-node path at the end of a round in the tree exploration game is only a partial play of the corresponding evaluation game, the truthifier simply outputs  $0^{\mathcal{M}}$  in the current evaluation game until a new game can be started.) In other words, when the truthifier adds a node  $v$  as the child of  $u$ , it can “replay” the path from the root to  $v$  using the moves  $m_1, m_2, \dots, m_d$  on the path, and the oblivious falsifier will choose the moves  $n_1, n_2, \dots, n_d$  as response. Therefore the truthifier can simulate the winning strategy for the tree exploration game by sequentially playing the evaluation game  $\ell$  times and beating the falsifier in at least one of the games.

We now describe in more detail the translation of an  $\mathcal{L}$ -term universal winning strategy in the tree exploration game into  $\mathcal{L}$ -strategies (with ancillary information) for the evaluation game. Consider the strategy  $\tau^{\text{tree}} = \langle p_1, r_1, p_2, r_2, \dots, p_\ell, r_\ell \rangle$ , where  $\ell$  is a constant. Recall that the location of each play of the truthifier in the tree exploration game is fixed, and that  $r_1, \dots, r_\ell$  describe the nodes to which a new child is added in each play. For each  $i \in [\ell]$ , we define an  $\mathcal{L}$ -term strategy for the evaluation game  $\tau_i^{\text{eval}}$  as follows:

- Let  $r_i$  be the node of the game tree that is extended during the  $i$ -th play of the tree exploration game. Suppose this node is at the  $d_i$ -th level of the tree, and let  $p_{i_1}, \dots, p_{i_{d_i}}$  be the  $\mathcal{L}$ -terms corresponding to the moves of the truthifier in the root-to- $r_i$  path, including the current move.
- Define the following  $\mathcal{L}$ -strategy  $\tau_i^{\text{eval}}$  of the evaluation game: (1) parse the ancillary information as a sequence  $\Gamma$  of transcripts derived from playing strategies  $\tau_1^{\text{eval}}, \dots, \tau_{i-1}^{\text{eval}}$  with the ancillary information described in the statement of the theorem; (2) in the  $j$ -th step (during the  $i$ -th evaluation game), where  $j \in [k]$ , if  $j \leq d_i$  play according to  $p_{i_j}$  using that all plays from previous rounds of the tree exploration game are available in the transcript  $\Gamma$ . Otherwise, choose  $0^{\mathcal{M}}$  (i.e., the  $j$ -th term defining the strategy is the constant term 0).

From the discussion above, the correctness of the translation is clear: if the strategies  $\tau_1^{\text{eval}}, \tau_2^{\text{eval}}, \dots, \tau_\ell^{\text{eval}}$  cannot beat a fixed falsifier strategy  $\tau^f$  in  $\ell$  sequential plays of the evaluation game, we can use the oblivious strategy defined by  $\tau^f$  in the tree exploration game to show that the truthifier does not win the tree exploration game withing  $\ell$  moves.  $\square$

*Remark 4.21.* Instead of viewing Theorem 4.20 as a special case of the game-theoretic witnessing theorem that employs the tree exploration game (Theorem 4.1), we can also establish the result in a more direct way using a technique known as the *no-counterexample interpretation*. We present a self-contained proof in Appendix C.

## 5 Warm-up: Krajíček’s Technique and the Pich-Santhanam Result

In this section, we provide a detailed exposition of the unprovability result from Santhanam and Pich [PS21], which relies on a technique introduced by Krajíček [Kra11] and further investigated by Pich [Pic15a]. Their result (intuitively) means that strong average-case circuit lower bounds against co-nondeterministic

circuits are not provable in  $\mathsf{TPV}$ . Concretely, for every  $L \in \mathsf{NTIME}[2^{n^{o(1)}}]$ ,  $\delta \in (0, 1) \cap \mathbb{Q}$ , and  $n_0 \in \mathbb{N}$ ,  $\mathsf{TPV}$  cannot prove that:

For every  $n > n_0$  and every co-nondeterministic circuit  $C : \{0, 1\}^n \rightarrow \{0, 1\}^1$  of size  $2^{n^\delta}$ ,  
 $C(x) = L(x)$  on at most  $\frac{1}{2} + \frac{1}{2^{n^\delta}}$  fraction of  $x \in \{0, 1\}^n$ .

Since our unprovability results are obtained by extending the original ideas of Pich and Santhanam [PS21] and Krajíček [Kra11] in combination with our new witnessing theorem, this section might be particularly helpful for a reader that is unfamiliar with these methods.

## 5.1 Formalization of Complexity Lower Bounds

While the unprovability result of [PS21] is robust to some details of the formalization, we will make a few comments here about the way it is done. First, we can represent any natural number  $a \in \mathbb{N}$  by an  $\mathcal{L}(\mathsf{PV})$ -term, e.g.,  $a = 1 + 1 + \dots + 1$ , where  $+$ :  $\mathbb{N}^2 \rightarrow \mathbb{N}$  is the  $\mathcal{L}(\mathsf{PV})$  function symbol for addition, and 1 is a constant symbol in  $\mathcal{L}(\mathsf{PV})$ . From this, we can introduce representations for other finite objects. For instance, a natural number can represent the code of a Turing machine  $M$ , while a pair of natural numbers can represent a rational number  $\delta \in \mathbb{Q}$ . In some cases, we will quantify over all such objects in the meta-language, e.g., if  $M$  is a Turing machine (in the usual sense), then we can consider a  $\mathcal{L}(\mathsf{PV})$ -sentence  $\phi_M$  that refers to machine  $M$  via its representation as a natural number.

For a nondeterministic Turing machine  $M$ , a constant  $n_0 \in \mathbb{N}$ , and functions  $s, m: \mathbb{N} \rightarrow \mathbb{N}$ , we write  $\mathsf{LB}(M, s, m, n_0)$  to denote an  $\mathcal{L}(\mathsf{PV})$ -sentence stating that, for every input length  $n \geq n_0$  and for every co-nondeterministic circuit  $D_n(x, z)$  of size  $\leq s(n)$ , there are at least  $m = m(n)$  distinct input strings  $x^1, \dots, x^m \in \{0, 1\}^n$  such that  $M(x^i) \neq D_n(x^i)$  for each  $1 \leq i \leq m$ .<sup>26</sup> A bit more formally, this sentence can be expressed in  $\mathcal{L}(\mathsf{PV})$  in the following way, where we assume that  $M$  on input length  $n$  runs in time  $\leq t(n)$  for some efficiently computable time bound  $t(n) \leq N(n) = 2^n$  and that  $s(n)$  and  $m(n)$  are efficiently computable and bounded by  $N = 2^n$ :

$$\begin{aligned} \mathsf{LB}(M, s, m, n_0) &\triangleq \forall v \forall N = |v| \forall n = |N| \text{ such that } n \geq n_0 \quad (\text{in other words, } n \in \mathsf{LogLog}) \\ &\quad \forall \text{ co-nondet. circuit } D_n \text{ of size } \leq s(n) \\ &\quad \exists m = m(n) \text{ distinct } n\text{-bit strings } x^1, \dots, x^m \text{ s.t. } \mathsf{Error}_{M, D_n}(x^i) \text{ for all } i \in [m], \end{aligned}$$

where we let  $\mathsf{Error}_{M, D_n}(x)$  denote the following  $\mathcal{L}(\mathsf{PV})$ -formula:

$$\mathsf{Error}_{M, D_n}(x) \equiv \left[ \exists y \exists z M(x, y) = 1 \wedge D_n(x, z) = 0 \right] \vee \left[ \forall y' M(x, y') = 0 \wedge \forall z' D_n(x, z') = 1 \right],$$

with the length of  $y, y'$  and  $z, z'$  bounded by the running time of  $M$  and the size of  $D_n$ , respectively.

The definition above can be made formal by the use of explicit  $\mathcal{L}(\mathsf{PV})$ -function symbols that evaluate circuits and machines on a given input and that perform other necessary checks, e.g., deciding when a given object represents a circuit of size at most  $s(n)$ . All this can be done without increasing the quantifier complexity of the resulting sentence, since  $n \in \mathsf{LogLog}$  and polynomial-time computations over  $N = 2^n$ -bit strings are feasible. For the same reason, the quantification over  $i \in [m]$  does not increase quantifier complexity, using that  $m(n) \leq N$ . Indeed, in the sentence it is enough to existentially quantify over  $m(n)$  strings  $x^i$  and over  $m(n)$  strings  $y^i, z^i$  followed by a universal quantification over  $m(n)$  strings  $y^i, z^i$ , and

<sup>26</sup>Here and throughout the exposition, we use that  $M(x) = 1$  if and only if there exists  $y$  such that  $M(x, y) = 1$  (as  $M$  computes nondeterministically), while  $D_n(x) = 1$  if and only if for every  $z$  we have  $D_n(x, z) = 1$  (since  $D$  is a co-nondeterministic circuit).



the remaining error conditions can be expressed using a single  $\mathcal{L}(\text{PV})$ -function symbol that gets as input the encoding of each collection of strings (formally, each family of  $m$  strings is a single object, and the strings are decoded from it). Overall, we get that  $\text{LB}(M, s, m, n_0)$  is a  $\forall\Sigma_2^b\text{-}\mathcal{L}(\text{PV})$  sentence.

**Theorem 5.1** ( $\text{T}_{\text{PV}}$  doesn't prove strong a.e. average-case co-nondeterministic lower bounds for NP). *For every  $n_0 \in \mathbb{N}$  and  $\delta \in \mathbb{Q} \cap (0, 1)$ , if  $M$  is a nondeterministic machine whose running time is bounded by some constructive function  $t(n) = 2^{n^{o(1)}}$ , then*

$$\text{T}_{\text{PV}} \not\vdash \text{LB}(M, s, m, n_0),$$

where  $s(n) = 2^{n^\delta}$  and  $m(n) = 2^n/2 - 2^n/2^{n^\delta}$ .

In particular, for every language  $L \in \text{NP}$  and  $\delta > 0$  it is consistent with  $\text{T}_{\text{PV}}$  that there are infinitely many input lengths  $n$  and a co-nondeterministic circuit  $D_n$  of size  $\leq 2^{n^\delta}$  such that

$$\Pr_{x \sim \{0,1\}^n} [L(x) = D_n(x)] \geq 1/2 + 2^{-n^\delta}.$$

A strengthening of Theorem 5.1 is discussed in Section 5.3.

## 5.2 Proof of Theorem 5.1

Let  $n_0, \delta, M, t(n), s(n)$ , and  $m(n)$  be as in the statement of Theorem 5.1. Arguing as in [PS21], we assume towards a contradiction that

$$\text{T}_{\text{PV}} \vdash \text{LB}(M, s, m, n_0).$$

Let  $L \subseteq \{0,1\}^*$  be the language defined by  $M$ . We argue as follows.

- (i) From the provability of this almost-everywhere average-case lower bound against co-nondeterministic circuits, it follows by the soundness of  $\text{T}_{\text{PV}}$  that (in the standard model) for every sequence  $\{E_n\}_{n \geq 1}$  of *deterministic* circuits  $E_n$  of size  $\leq 2^{n^\delta}$ , if  $n \geq n_0$  then

$$\Pr_{x \sim \{0,1\}^n} [L(x) = E_n(x)] \leq 1/2 + 2^{-n^\delta}.$$

- (ii) From the provability of the sentence  $\text{LB}(M, s, m, n_0)$  it trivially follows that  $\text{T}_{\text{PV}}$  proves a sentence  $\text{LB}_{\text{wst}}(M, s, n_0)$  which states a *worst-case* lower bound for  $M$  against co-nondeterministic circuits of the same size. We then show that the provability of  $\text{LB}_{\text{wst}}(M, s, n_0)$  in  $\text{T}_{\text{PV}}$  implies that, in the standard model, for every fixed  $k \geq 1$  and for every large enough  $n$ , there is a deterministic circuit  $A$  defined over  $n^k$  input variables and of size  $2^{O(n)}$  such that

$$\Pr_{w \sim \{0,1\}^{n^k}} [L(w) = A(w)] \geq 1/2 + 2^{-O(n)}.$$

Taking  $k > 1/\delta$  contradicts Item (i) above.

Note that the only remaining step is to show that:

- (\*) The provability of a worst-case lower bound against *co-nondeterministic* circuits allows us to non-trivially approximate  $L$  using *deterministic* circuits of bounded size.

Before proceeding with the proof of this result, we describe the aforementioned worst-case lower bound sentence in a convenient way.

$$\begin{aligned} \text{LB}_{\text{wst}}(M, s, n_0) \equiv & \forall n \in \text{LogLog with } n \geq n_0, \forall \text{ co-nondet. circuit } D \text{ of size } \leq s(n) \\ & \exists x \in \{0, 1\}^n \exists y \in \{0, 1\}^{t(n)} \exists z \in \{0, 1\}^{s(n)} \text{ such that Error}(x, y, z), \end{aligned}$$

where here  $\text{Error}(x, y, z)$  denotes the following  $\mathcal{L}(\text{PV})$ -formula:

$$\text{Error}(x, y, z) \equiv \left[ M(x, y) = 1 \wedge D(x, z) = 0 \right] \vee \left[ \forall y' M(x, y') = 0 \wedge \forall z' D(x, z') = 1 \right], \quad (5)$$

where the lengths of  $y'$  and  $z'$  are bounded as before. Observe that  $\text{LB}_{\text{wst}}(M, s, n_0)$  is also a  $\forall\Sigma_2^b\text{-}\mathcal{L}(\text{PV})$  sentence.

It is easy to see that, under any reasonable formalization, if  $m(n) \geq 1$  then  $\text{T}_{\text{PV}}$  derives the worst-case lower bound sentence  $\text{LB}_{\text{wst}}(M, s, n_0)$  from the average-case lower bound sentence  $\text{LB}(M, s, m, n_0)$ . Consequently, it is sufficient for us to prove the following lemma, which formalizes statement  $(\star)$ .

**Lemma 5.2** (Non-trivial correlation from the provability of a worst-case lower bound). *Let  $n_0 \in \mathbb{N}$ ,  $\delta \in \mathbb{Q} \cap (0, 1)$ ,  $M$  be a nondeterministic machine whose running time is bounded by some constructive function  $t(n) = 2^{n^{o(1)}}$ , and  $s(n) = 2^{n^\delta}$ . If*

$$\text{T}_{\text{PV}} \vdash \text{LB}_{\text{wst}}(M, s, n_0),$$

*then for every  $k \geq 1$  and sufficiently large  $n$ , there is a deterministic circuit  $B : \{0, 1\}^{n^k} \rightarrow \{0, 1\}$  of size  $2^{O(n)}$  such that*

$$\Pr_{w \sim \{0, 1\}^{n^k}} [L(w) = B(w)] \geq 1/2 + 2^{-O(n)},$$

*where  $L$  is the language decided by  $M$ .*

Lemma 5.2 and its proof have appeared in [Kra11, Pic15a]. We provide next a detailed exposition of this technique.

### 5.2.1 A simpler case: $\ell = 1$ in the KPT student-teacher protocol

Note that we can apply Theorem 3.11 to sentence  $\text{LB}_{\text{wst}}$  and theory  $\text{T}_{\text{PV}}$ , since  $\text{T}_{\text{PV}}$  is a universal theory and it is not difficult to see that  $\text{LB}_{\text{wst}}$  can be written in the form required by Theorem 3.11. Since  $\mathcal{L}(\text{PV})$ -terms correspond to polynomial-time computable functions, the corresponding student computes in time polynomial in the length of its input. Using that  $n \in \text{LogLog}$  in sentence  $\text{LB}_{\text{wst}}$ , we obtain uniform algorithms  $f_1, \dots, f_\ell$  that compute in time  $2^{O(n)}$  and satisfy the conclusion of Theorem 3.11. Note that we cannot control the constant  $\ell$ . In this section, we discuss the simpler case when we get  $\ell = 1$  in the application of Theorem 3.11 to  $\text{T}_{\text{PV}} \vdash \text{LB}_{\text{wst}}(M, s, n_0)$ .

Omitting auxiliary variables in the input to  $f_1$  and highlighting the relevant parameters,<sup>27</sup>  $f_1(n, D)$  receives  $n$  and an arbitrary co-nondeterministic circuit  $D$  of size  $\leq s(n) = 2^{n^\delta}$ , and outputs a triple  $(x, y, z)$  such that  $\text{Error}(x, y, z)$  holds. Assuming that  $n \geq n_0$  and  $\ell = 1$  in the KPT Witnessing, it follows that this

<sup>27</sup>The condition  $n \in \text{LogLog}$  means that  $n$  is the length of  $N$ , while  $N$  is the length of some universally quantified variable  $v$ . For this reason, formally,  $v$  is an input to  $f_1$ , and not  $n$ . However, over  $\mathbb{N}$  we will run  $f_1$  on a fixed input  $v$ , such as  $1^N$ , where  $N = 2^n$ . For this reason,  $n$  is the parameter that controls the length of the remaining inputs.

triple witnesses that  $D(x) \neq M(x)$  over the standard model  $\mathbb{N}$ . In other words,  $x \in \{0, 1\}^n$ ,  $y \in \{0, 1\}^{t(n)}$  for  $t(n) = 2^{n^{o(1)}}$ ,  $z \in \{0, 1\}^{s(n)}$ , and the following holds:

$$\left[ M(x, y) = 1 \wedge D(x, z) = 0 \right] \vee \left[ \forall y' M(x, y') = 0 \wedge \forall z' D(x, z') = 1 \right].$$

To prove Lemma 5.2 when  $\ell = 1$ , let  $k \geq 1$ , and assume that  $n$  is sufficiently large. We will use  $f_1$  to construct a deterministic circuit  $B$  defined over  $n^k$  input variables and of size  $2^{O(n)}$  such that

$$\Pr_{w \sim \{0, 1\}^{n^k}} [L(w) = B(w)] \geq 1/2 + 2^{-O(n)},$$

where  $L$  is the language computed by the nondeterministic machine  $M$ . As a key point, note that we can invoke  $f_1(n, D)$  on any co-nondeterministic circuit  $D(x, \cdot)$  over  $n$  input variables and of size  $\leq 2^{n^\delta}$ . In order to construct the deterministic circuit  $B$ , we will also use that  $f_1(n, D) = (x, y, z)$  computes in time  $2^{O(n)}$  and therefore can be simulated by a circuit of size  $2^{O(n)}$ . For the  $\ell = 1$  case, we will not need the output strings  $y$  and  $z$  during the construction of  $A$ . From now on, we simplify notation and write  $f_1(D) = x$  to denote the relevant input and output of  $f_1$  for this case.

Let  $C_{n^k}(w, z)$  be a Boolean circuit that computes as  $M$  on inputs of length  $n^k$ , where  $z$  corresponds to the nondeterministic input. Since  $M$  runs in time at most  $2^{n^{o(1)}}$ ,  $C_{n^k}$  has size at most  $2^{n^\delta}$  when  $n$  is sufficiently large. We partition its first input as  $w = x \| w'$ , where  $|x| = n$  and  $|w'| = n^k - n$ . Now for a fixed string  $w' \in \{0, 1\}^{n^k - n}$ , consider the circuit  $D_{w'}(x, z) \triangleq \neg C_{n^k}(x \| w', z)$ . Viewing  $D_{w'}$  as a co-nondeterministic circuit, we get that

$$D_{w'}(x) = 1 \iff \forall z D_{w'}(x, z) = 1 \iff \forall z C_{n^k}(x \| w', z) = 0 \iff x \| w' \notin L(M). \quad (6)$$

Intuitively, if we “learn” the output bit of  $D_{w'}(x)$  for some pair  $(w', x)$ , we also “learn” if the input string  $x \| w'$  is in  $L(M)$ . As a consequence, the collection  $\{D_{w'}(x)\}_{w'}$  of co-nondeterministic circuits  $D_{w'}$  (defined over  $n$  input bits) captures the computation of  $L$  on inputs of length  $n^k$ .

As each  $D_{w'}$  has size at most  $2^{n^\delta}$ , we can invoke  $f_1$  on them. Since  $f_1(D_{w'})$  finds mistakes with respect to the nondeterministic computation of  $M$ , we know that  $D_{w'}(x^*) \neq M(x^*)$  for  $x^* \triangleq f_1(D_{w'})$ . Since there are only  $2^n$  possibilities for the output of  $f_1$ , the following holds.

**Fact 5.3.** *There is a string  $x^* \in \{0, 1\}^n$  such that*

$$\Pr_{w' \in \{0, 1\}^{n^k - n}} [f_1(D_{w'}) = x^*] \geq 2^{-n}.$$

Following our informal discussion from above, from the knowledge that  $f_1(D_{w'}) = x^*$  (and of a bit encoding if  $x^* \in L(M)$ ) we “learn” how to compute  $L(M)$  on the input  $w = x^* \| w'$ . Since this will happen with non-trivial probability over a random choice of  $w'$  and  $x^*$ , this can be used to non-trivially approximate  $L(M)$  over input length  $n^k$ .

Formally, let  $b^* \in \{0, 1\}$  be 1 if and only if  $M(x^*) = 1$ , i.e., if  $x^* \in L(M)$ . The string  $x^*$  and the bit  $b^*$  will be stored as non-uniform advice in the deterministic circuit  $B$  that we show to be correlated with  $L(M)$  on input length  $n^k$ . First, consider the following randomized circuit  $B'$ :

Note that  $B'$  can be computed with  $2^{O(n)}$  gates, since  $f_1$  runs in time  $2^{O(n)}$ . Next, we show that the randomized circuit  $B'$  non-trivially correlates with  $L(M)$  on inputs of length  $n^k$ . After that, fixing the random bit  $b$  in  $B'$  yields the desired deterministic circuit  $A$ .

**Input :** The input  $w \in \{0, 1\}^{n^k}$  and a random  $r \in \{0, 1\}$

**Advice:**  $x^* \in \{0, 1\}^n$  and  $b^* \in \{0, 1\}$

1 Let  $w = x \| w'$  and compute  $f_1(D_{w'})$ ;

// Note that we can construct a description of  $D_{w'}$  from  $w$  and  $C_{n^k}$ .

2 If  $f_1(D_{w'}) \neq x^*$ , **return**  $r$ ;

3 If  $x \neq x^*$ , **return**  $r$ ;

4 Otherwise, **return**  $b^*$ ;

**Algorithm 1:** Randomized Circuit  $B'$  for  $L(M)$

**Fact 5.4.** If  $B'$  reaches Line 4 on an input string  $w$ , then  $B'(w, b) = M(w)$ , i.e.,  $B'$  correctly decides  $L(M)$  on input  $w$ .

*Proof.* Under the assumption that  $B'$  reaches Line 4 on an input string  $w$ , it follows that  $w = x^* \| w'$  and  $f_1(D_{w'}) = x^*$ . Moreover, observe that the random bit  $b$  does not affect the output of  $B'$  in this case. We have

$$\begin{aligned}
 B'(w, r) = 1 &\iff b^* = 1 \\
 &\iff M(x^*) = 1 && \text{(by the definition of } b^*) \\
 &\iff D_{w'}(x^*) = 0 && \text{(using that } f_1 \text{ finds a mistake)} \\
 &\iff w = x^* \| w' \in L(M) && \text{(by Equation 6)} \\
 &\iff M(w) = 1. && \text{(since } M \text{ computes } L(M))
 \end{aligned}$$

In other words,  $B'(w, r) = M(w)$ .  $\square$

In addition, by Fact 5.3,

$$p \triangleq \Pr_{r,w}[B' \text{ reaches Line 4}] = \Pr_{x,w'}[x = x^* \wedge f_1(D_{w'}) = x^*] \geq 2^{-n} \cdot 2^{-n}.$$

On the other hand, when  $B'$  does not reach Line 4 it outputs a random bit that is independent of the input string  $w$ . Therefore, using Fact 5.4 and the lower bound on  $p$ ,

$$\Pr_{r,w}[B'(w, r) = M(w)] \geq p \cdot 1 + (1 - p) \cdot 1/2 = 1/2 + p/2 \geq 1/2 + 2^{-2n+1} = 1/2 + 2^{-O(n)}.$$

Fixing the random bit  $r$  in the best way maintains this advantage and completes the proof of Lemma 5.2 when  $\ell = 1$ .

*Remark 5.5.* The same argument can be used to approximate any nondeterministic circuit of size  $2^{n^{k\delta}}$  defined over  $n^k$  bits by a deterministic circuit of size  $2^{O(n)}$ , instead of just for  $L(M) \cap \{0, 1\}^{n^k}$ . In other words, by connecting  $D_{w'}$  to the computation of the appropriate co-nondeterministic circuit, “learning” output bits of  $D_{w'}$  via  $f_1$  translates into a non-trivial approximation (using exactly the same strategy). This will also hold when analysing the case  $\ell > 1$ . In particular, from  $\text{TPV} \vdash \text{LB}_{\text{wst}}(M, s, n_0)$  we are able to non-trivially approximate any language in  $\text{NSIZE}[2^{n^{o(1)}}]$  and not just  $L(M)$ .

### 5.2.2 The case $\ell = 2$ via the Nisan-Wigderson generator

In this section, we consider the case where the disjunction obtained from KPT Witnessing (Theorem 3.11) has size  $\ell = 2$ . This essentially covers all difficulties in the general case. Before handling  $\ell = 2$ , it is instructive to highlight some key points of the proof when  $\ell = 1$ :

- (i) We implicitly relied on the ability of *certifying* when an input  $x$  is a mistake. More precisely, when  $\ell = 1$ , if  $f_1(D_{w'}) = (x, y, z)$ , we have the *guarantee* that  $D_{w'}(x) \neq M(x)$ . This is because there is a *single* round in the corresponding Student-Teacher protocol.
- (ii) By an averaging argument, we fixed a good string  $x^* \in \{0, 1\}^n$  (Fact 5.3), which eventually allowed us to compute  $M(x^*w')$  on a non-trivial fraction of  $w'$ , by storing  $x^*$  and the corresponding bit  $b^* = M(x^*)$ .
- (iii) This was accomplished by considering a family  $\{D_{w'}(x)\}_{w'}$  of co-nondeterministic circuits over  $n$ -bit inputs that compute according to a circuit defined over input length  $n^k$  that is related to the language we would like to approximate.
- (iv) On an input  $w \in \{0, 1\}^{n^k}$  with  $w = x\|w'$  for which the witnessing provided by  $f_1(D_{w'})$  was inconsistent with the actual input part  $x$  (we can easily detect this), we output a random bit.

Note that this approach no longer works when  $\ell > 1$ : the first term obtained from KPT Witnessing might not succeed in finding a mistake. For this reason, we cannot assume in Item (i) that if  $f_1(D_{w'}) = (x, y, z)$  then  $D_{w'}(x) \neq M(x)$ .

Let  $f_1$  be the first term in the KPT disjunction when  $\ell > 1$ . Note that we can still fix a popular *candidate* mistake  $x^* \in \{0, 1\}^n$ , as in Fact 5.3. Recall that  $f_1(D_{w'}) = (x_{w'}, y_{w'}, z_{w'})$  (we did not have to use  $y_{w'}$  and  $z_{w'}$  in the argument for  $\ell = 1$ ). We can check whether  $x_{w'} = x^*$ , as before, and we would like to use  $y_{w'}$  and  $z_{w'}$  together with some hard-coded information to decide if  $x^* = x_{w'}$  is indeed a mistake for  $D_{w'}$ . While both  $y_{w'}, z_{w'} \in \{0, 1\}^{\leq 2^n}$ , there are  $2^{\Omega(n^k)}$  possible strings  $w'$ . Unfortunately, it is unclear how to store enough information in the non-uniform circuit  $B'$  to certify that a mistake has been found by  $f_1$  while maintaining a circuit size bound of  $2^{O(n)}$ .

To reduce the amount of advice needed in  $B'$  and address this difficulty, the solution [Kra11, Pic15a] is to employ a more sophisticated family  $\{D_{w'}\}_{w'}$  of circuits constructed via the Nisan-Wigderson generator [NW94].

For a nondeterministic machine  $M$  that decides a language  $L(M)$ , we use the notation  $\{\text{NW}_{\overline{L(M)}}(w)\}_w$  to denote the collection of functions obtained from the Nisan-Wigderson generator when instantiated with the Boolean function  $h$  that corresponds to the negation of  $L(M)$  over inputs of length  $n^{c/2}$ .

**Fact 5.6.** *Let  $M$  be a nondeterministic machine that runs in time  $2^{n^{o(1)}}$  on inputs of length  $m$ . For any constant  $c \geq 1$  and every large enough  $n$ , each function in  $\{\text{NW}_{\overline{L(M)}}(w)\}_w$  can be computed by a co-nondeterministic circuit  $D_w(x)$  of size at most  $2^{n^\delta}$ .*

**The case  $\ell = 1$  via the NW generator.** Before handling the case  $\ell = 2$ , we sketch the proof of the case  $\ell = 1$  using the collection  $\{D_w\}_{w \in \{0, 1\}^{n^c}}$  obtained from the nondeterministic machine  $M$  and the NW generator, with parameters as above.

Consider the function  $f_1(D_w) = (x, y, z)$  obtained by applying Theorem 3.11, and assume that  $\ell = 1$ . Again, we will not inspect  $y$  and  $z$  when  $\ell = 1$ . Recall that  $f_1(D_w)$  computes in time  $2^{O(n)}$ . We show how to decide  $L(M)$  on inputs of length  $n^{c/2}$  by a deterministic circuit of size  $2^{O(n)}$  that agrees with  $L(M)$  with probability  $\geq 1/2 + 2^{-O(n)}$  over a uniformly random input string.

Similarly to Fact 5.3, by a standard averaging argument we can establish the following fact.

**Fact 5.7.** *There is a string  $x^* \in \{0, 1\}^n$  such that*

$$\Pr_{w \in \{0, 1\}^{n^c}} [f_1(D_w) = x^*] \geq 2^{-n}.$$

Recall that  $J_{x^*}$  denotes the subset of  $[n^c]$  of size  $n^{c/2}$  corresponding to the  $x^*$ -row of the design in our NW generator; for  $a \in \{0, 1\}^{n^c - n^{c/2}}$  and  $u \in \{0, 1\}^{n^{c/2}}$ ,  $r_x(a, u)$  denotes the “concatenated” string  $a \cup u$  obtained by viewing  $a \in \{0, 1\}^{[n^c] \setminus J_{x^*}}$  and  $u \in \{0, 1\}^{J_{x^*}}$ . By another averaging argument, we get the following consequence.

**Fact 5.8.** *There is a string  $a \in \{0, 1\}^{[n^c] \setminus J_{x^*}}$  of length  $n^c - n^{c/2}$  such that*

$$\Pr_{\substack{u \sim \{0, 1\}^{J_{x^*}} \\ w \triangleq a \cup u}} [f_1(D_w) = x^*] \geq 2^{-n}.$$

We can view  $D_w = \text{NW}_{\overline{L(M)}}(w)$  as a co-nondeterministic circuit for computing  $\overline{L(M)}$  over inputs of length  $n^{c/2}$  derived from the seed  $w$ :

$$D_w(x) = 1 \iff w|_{J_{x^*}} \in \overline{L(M)}.$$

Given the previous discussion, we are interested in seeds  $w \in \{0, 1\}^{n^{2c}}$  of the form  $w = a \cup u$ , where  $a \in \{0, 1\}^{[n^c] \setminus J_{x^*}}$  is fixed,  $u \in \{0, 1\}^{J_{x^*}}$ , and  $f_1(D_w) = x^*$ . We know that a non-trivial fraction of strings  $u$  will satisfy this condition. Since  $f_1$  witnesses mistakes with respect to  $L(M)$  over inputs of length  $n$  (note that  $D_w$  is a conondeterministic circuit over  $n$ -bit inputs), whenever  $f_1(D_w) = x^*$  we are guaranteed that

$$D_w(x^*) = 1 \iff M(x^*) = 0,$$

which implies that  $M(x^*) = 0$  if and only if  $w|_{J_{x^*}} \notin L(M)$ . Now  $x^*$  is fixed, so the equality  $M(x^*) = 0$  does not depend on other conditions. For instance, if  $M(x^*) = 0$ , we can conclude that on any input string  $u \in \{0, 1\}^{n^{c/2}}$ , if for  $w = a \cup u$  we have  $f_1(D_w) = x^*$ , then  $u = w|_{J_{x^*}}$  is not in  $L(M)$ . Consequently, this allows us to correctly compute  $L(M)$  on any such input  $u \sim \{0, 1\}^{n^{c/2}}$ , which constitute a non-trivial fraction of inputs. Moreover, we can check whether an input  $u$  satisfies  $f_1(D_w) = x^*$  using a deterministic circuit of size  $2^{O(n)}$ .

Formally, consider the fixed strings  $x^* \in \{0, 1\}^n$  and  $a \in \{0, 1\}^{[n^c] \setminus J_{x^*}}$  from above, and let  $b^* \triangleq M(x^*) \in \{0, 1\}$ . We hardcode  $x^*$ ,  $a$ , and  $b^*$  in the randomised circuit  $B(u)$  described below:

**Input :** The input  $u \in \{0, 1\}^{n^{c/2}}$  and a random  $r \in \{0, 1\}$   
**Advice:**  $x^* \in \{0, 1\}^n$  and  $b^* \in \{0, 1\}$

- 1 Let  $w = r_{x^*}(a, u)$ ;
- 2 Let  $x = f_1(D_w)$ ;
- 3 If  $x \neq x^*$ , output the random bit  $r$ ;
- 4 Otherwise, **return**  $b^*$ ;

**Algorithm 2:** Randomised Circuit  $B$  for  $L(M)$

Given the aforementioned discussion, it is easy to see that

$$\Pr_{u, r} [B(u, r) = M(u)] \geq p \cdot 1 + (1 - p) \cdot 1/2 = 1/2 + p/2 \geq 1/2 + 2^{-n+1},$$

where  $p$  is the probability in the LHS of Fact 5.8. Consequently, by an averaging argument over the random bit  $r$ , there is a deterministic circuit of size  $2^{O(n)}$  that computes  $L(M)$  on inputs of length  $n^{c/2}$  with the same advantage.



**The case  $\ell = 2$  via the NW generator.** Recall that

$$\text{Error}(x, y, z) \equiv \left[ M(x, y) = 1 \wedge D(x, z) = 0 \right] \vee \left[ \forall y' M(x, y') = 0 \wedge \forall z' D(x, z') = 1 \right].$$

We now have a function  $f_1(D) = (x, y, z)$  that attempts to produce a triple  $(x, y, z)$  satisfying  $\text{Error}(x, y, z)$ , and a function  $f_2(D, y', z')$  which given a pair  $y', z'$  for which

$$\left[ M(x, y) = 0 \vee D(x, z) = 1 \right] \wedge \left[ M(x, y') = 1 \vee D(x, z') = 0 \right] \quad (7)$$

is able to produce an input  $x'$  such that  $D(x') \neq M(x')$ .

Again, we consider the family  $\{D_w\}_{w \in \{0,1\}^{n^c}}$  of conondeterministic circuits  $D_w$  of size  $\leq 2^{n^\delta}$  that compute  $\text{NW}_{\overline{L(M)}}(w): \{0,1\}^n \rightarrow \{0,1\}$  for a fixed seed  $w$ , with parameters as described above. In particular, this generator is instantiated with respect to the Boolean function  $h$  corresponding to  $\overline{L(M)}$  over inputs of length  $n^{c/2}$ , for a fixed but arbitrarily large constant  $c \geq 1$ .

By an averaging argument, the following claim holds.

**Fact 5.9.** *There is a string  $x_1 \in \{0,1\}^n$  such that*

$$\Pr_{w \in \{0,1\}^{n^c}} [f_1(D_w) = x_1] \geq 2^{-n}.$$

Fix this  $x_1$ . We define the sets  $S_{x_1}^{\text{mist}} \subseteq S_{x_1} \subseteq \{0,1\}^{n^c}$  as follows:

$$\begin{aligned} S_{x_1} &\triangleq \left\{ w \in \{0,1\}^{n^c} \mid f_1(D_w) = x_1 \right\}, \\ S_{x_1}^{\text{mist}} &\triangleq \left\{ w \in S_{x_1} \mid D_w(x_1) \neq M(x_1) \right\}, \end{aligned}$$

and consider the density of  $S_{x_1}^{\text{mist}}$  with respect to its superset  $S_{x_1}$ .

**Case 1.**  $|S_{x_1}^{\text{mist}}| > (2/3) \cdot |S_{x_1}|$ . We can essentially proceed as in the case of  $\ell = 1$ , with the exception that one needs to be careful when invoking an analogue of Fact 5.8. This is because fixing a string  $a \in \{0,1\}^{[n^c] \setminus J_{x_1}}$  might keep the density of  $S_{x_1}$  at least  $2^{-n}$  but could significantly decrease the relative density of the set  $S_{x_1}^{\text{mist}}$  after the restriction.

To handle this, we introduce the following notation. For  $m \geq 1$ , a set  $S \subseteq \{0,1\}^{[m]}$ , and a string  $a \in \{0,1\}^I$ , where  $I \subseteq [m]$ , we define the *restriction of  $S$  with respect to  $a$*  as the set

$$S \upharpoonright_a \triangleq \{w \in S \mid w|_I = a\}.$$

Under the assumption that  $|S_{x_1}^{\text{mist}}| > (2/3) \cdot |S_{x_1}|$ , it is possible to show by a counting argument (see, e.g., Lemma E.1) that there exists a string  $a \in \{0,1\}^{[n^c] \setminus J_{x_1}}$  such that

$$p \triangleq \frac{|S_{x_1} \upharpoonright_a|}{2^{n^{c/2}}} \geq \frac{1}{n} \cdot 2^{-n} \quad \text{and} \quad \frac{|S_{x_1}^{\text{mist}} \upharpoonright_a|}{|S_{x_1} \upharpoonright_a|} \geq \frac{2}{3} - \frac{1}{n}. \quad (8)$$

While it is not clear how to decide in size  $2^{O(n)}$  if a string  $w \in S_{x_1}^{\text{mist}} \upharpoonright_a$ , we can check whether  $w \in S_{x_1} \upharpoonright_a$ . Since  $S_{x_1}^{\text{mist}}$  is dense in  $S_{x_1} \upharpoonright_a$ , this is enough to adapt the original strategy used for  $\ell = 1$ .

Formally, fix strings  $x_1 \in \{0,1\}^n$  and  $a \in \{0,1\}^{[n^c] \setminus J_{x_1}}$  as above, and let  $b_1 \triangleq M(x_1) \in \{0,1\}$ . We hardcode  $x_1, a$ , and  $b_1$  in the randomised circuit  $B_1(u, r)$  described below.

**Input :** The input  $u \in \{0, 1\}^{n^{c/2}}$  and a random  $r \in \{0, 1\}$

**Advice:**  $x_1 \in \{0, 1\}^n$ ,  $a \in \{0, 1\}^{n^c - n^{c/2}}$ , and  $b_1 \in \{0, 1\}$

- 1 Let  $w = r_{x_1}(a, u)$ ;
- 2 Let  $x = f_1(D_w)$ ;
- 3 If  $x \neq x_1$ , output the random bit  $r$ ;
- 4 Otherwise, output the fixed bit  $b_1 = M(x_1)$ ;

**Algorithm 3:** Randomized Circuit  $B_1$  for  $L(M)$  when  $|S_{x_1}^{\text{mist}}| > (2/3) \cdot |S_{x_1}|$ .

Clearly,  $B_1$  is computed by a randomised circuit of size  $2^{O(n)}$ . To analyse its success probability, first note that if  $u$  is such that  $w = r_{x_1}(a, u) \notin S_{x_1} \upharpoonright_a$ , then  $B_1(u) = M(u)$  with probability  $1/2$ . On the other hand, for those  $u$  such that  $w = r_{x_1}(a, u) \in S_{x_1} \upharpoonright_a$ , at least a  $2/3 - 1/n$  fraction of them are in  $S_{x_1}^{\text{mist}} \upharpoonright_a$ , in which case  $B_1(u)$  is correct. Since  $S_{x_1} \upharpoonright_a$  has density at least  $1/n \cdot 2^{-n}$ , it follows that

$$\Pr_{u,r}[B_1(u, r) = M(u)] = (1 - p) \cdot \frac{1}{2} + p \cdot \left(\frac{2}{3} - \frac{1}{n}\right) = \frac{1}{2} + p \cdot \left(\frac{1}{6} - \frac{1}{n}\right) = \frac{1}{2} + \Omega\left(\frac{2^{-n}}{n}\right),$$

which is  $1/2 + 2^{-O(n)}$ . Fixing the random bit  $r$  in the best way yields the desired deterministic circuit.

**Case 2.**  $|S_{x_1}^{\text{mist}}| < (2/3) \cdot |S_{x_1}|$ . In this case, the mistakes of at least a  $1/3$  fraction of the circuits  $D_w$  for  $w \in S_{x_1}$  must be witnessed by  $f_2$ . To make sure the output of  $f_2(D_w, y', z')$  is indeed a string  $x_2$  for which  $M(x_2) \neq D_w(x_2)$ , we must provide a pair  $y', z'$  such that

$$\left[ M(x_1, y_1) = 0 \vee D_w(x_1, z_1) = 1 \right] \wedge \left[ M(x_1, y') = 1 \vee D_w(x_1, z') = 0 \right], \quad (9)$$

where  $f_1(D_w) = (x_1, y_1, z_1)$ . We consider the Teacher that to each  $w \in S_{x_1} \setminus S_{x_1}^{\text{mist}}$  and corresponding  $(y_1, z_1)$  assign the lexicographic first pair  $(y'_w, z'_w)$  for which Equation (9) holds. Note that such a pair always exists, since in this case for  $x_1 = f_1(D_w)$  we have  $M(x_1) = D_w(x_1)$ .

By an averaging argument, the following claim holds.

**Fact 5.10.** *Under this fixed Teacher, there is a string  $x_2 \in \{0, 1\}^n$  such that the set*

$$S_{x_1, x_2} \triangleq \{w \in S_{x_1} \mid D_w(x_1) = M(x_1) \wedge f_2(D_w, y'_w, z'_w) = x_2\}$$

*has density at least  $(1/3) \cdot 2^{-2n}$  in  $\{0, 1\}^{n^c}$ .*

Note that, by construction, if  $w \in S_{x_1, x_2}$  then for  $x_2 = f_2(D_w, y'_w, z'_w)$  we have  $D_w(x_2) \neq M(x_2)$ . Note that  $x_1 \neq x_2$  because otherwise we have  $S_{x_1, x_2} = \emptyset$ .<sup>28</sup> Moreover, the set  $S_{x_1, x_2}$  has enough density for our purposes. However, for this to be useful we must verify that a given circuit  $D_w$  satisfies  $w \in S_{x_1, x_2}$  using a deterministic circuit of size  $2^{O(n)}$ .

By another averaging argument, we have the following result.

**Fact 5.11.** *There is a string  $a \in \{0, 1\}^{[n^c] \setminus J_{x_2}}$  such that*

$$\frac{|S_{x_1, x_2} \upharpoonright_a|}{2^{n^{c/2}}} \geq \frac{1}{3} \cdot 2^{-2n}.$$

<sup>28</sup> Assume it is not the case, there is a  $w \in S_{x_1, x_2}$  such that  $D_w(x_2) \neq M(x_2)$ . However, we know that  $D_w(x_1) = M(x_1)$  by the definition of  $S_{x_1, x_2}$ , which is impossible when  $x_1 = x_2$ .

Fix this string  $a \in \{0, 1\}^{[n^c] \setminus J_{x_2}}$  together with the strings  $x_1$  and  $x_2$ . We will assume that the following computation is possible in order to complete the proof, returning to it later on:

- ( $\nabla$ ) There is a deterministic circuit  $E(w)$  of size  $2^{O(n)}$  as follows: Given a  $w \in S_{x_1}$  of the form  $a \cup u$  such that  $D_w(x_1) = M(x_1)$ , it outputs the lexicographic first pair  $(y'_w, z'_w)$  for which Equation (9) holds, where  $(x_1, y_1, z_1) = f_1(D_w)$ .<sup>29</sup>

Consider strings  $x_1, x_2 \in \{0, 1\}^n$  and  $a \in \{0, 1\}^{[n^c] \setminus J_{x_2}}$  as above, and let  $b_2 \triangleq M(x_2) \in \{0, 1\}$ . We hardcode this information in the randomised circuit  $B_2(u)$  described below, which includes the circuit  $E(w)$  from ( $\nabla$ ) as a subroutine:

**Input :** The input  $u \in \{0, 1\}^{n^{c/2}}$  and a random  $r \in \{0, 1\}$   
**Advice:**  $x_1, x_2 \in \{0, 1\}^n$ ,  $a \in \{0, 1\}^{n^c - n^{c/2}}$ , and  $b_2 \in \{0, 1\}$

- 1 Let  $w = r_{x_2}(a, u)$ ;
- 2 Let  $(x, y_1, z_1) = f_1(D_w)$ ;
- 3 If  $x \neq x_1$ , output the random bit  $r$ ;
- 4 Let  $(y'_w, z'_w) = E(w)$ ;
- 5 If the tuple  $(x_1, y_1, z_1, y'_w, z'_w)$  satisfies Equation (9) and  $f_2(D_w, y'_w, z'_w) = x_2$ , output  $b_2$ ;
- 6 Otherwise output the random bit  $r$ .

**Algorithm 4:** Randomized Circuit  $B_2$  for  $L(M)$  when  $|S_{x_1}^{\text{mist}}| \leq (2/3) \cdot |S_{x_1}|$ .

Note that, under assumption ( $\nabla$ ),  $B_2$  can be computed by a randomised circuit of size  $2^{O(n)}$ . Moreover, it follows from our discussion and from the density of  $S_{x_1, x_2} \upharpoonright_a$  that

$$\Pr_{u, r}[B_2(u, r) = M(u)] \geq \frac{1}{2} + \Omega(2^{-2n}).$$

This yields a deterministic circuit with the same advantage.<sup>30</sup>

*Proof of ( $\nabla$ ).* We will now use the main property of the combinatorial design behind the Nisan-Wigderson generator: the sets  $J_{x_1}$  and  $J_{x_2}$  overlap in at most  $n$  coordinates. This will allow us to hardcode all relevant pairs  $(y'_w, z'_w)$  using circuit size  $2^{O(n)}$ .

To implement ( $\nabla$ ), we are given a string  $w = a \cup u$ , where  $a \in \{0, 1\}^{[n^c] \setminus J_{x_2}}$  is fixed and  $u \in \{0, 1\}^{J_{x_2}}$ , such that the following conditions hold:

- Let  $(x, y_1, z_1) = f_1(D_w)$ , then  $x = x_1$ .
- $D_w(x_1) = M(x_1)$ .

Our goal is to output the lexicographic first pair  $(y'_w, z'_w)$  such that:

$$\left[ M(x_1, y_1) = 0 \vee D_w(x_1, z_1) = 1 \right] \wedge \left[ M(x_1, y') = 1 \vee D_w(x_1, z') = 0 \right].$$

<sup>29</sup>Note that in this case  $f_2(D_w, y'_w, z'_w)$  outputs a mistake of  $D_w$ , since  $\text{Error}(x_1, y_1, z_1)$  does not hold and correct witnesses for this are provided.

<sup>30</sup>Note that we cannot really guarantee that  $D_w(x_1) = M(x_1)$  when invoking ( $\nabla$ ), since this cannot be easily decided in deterministic size  $2^{O(n)}$ . This means that more inputs  $u$  than those leading to strings  $w \in S_{x_1, x_2} \upharpoonright_a$  might reach Line 5 and be assigned output value  $b_2$ . Nevertheless,  $B_2$  will be correct on any such input  $u$ , by virtue of the two checks performed in Line 5. Put another way, the argument “covers” the inputs  $u$  leading to strings  $w \in S_{x_1, x_2} \upharpoonright_a$ .

Note that such pair must exist since we assume that  $D_w(x_1) = M(x_1)$ .

Recall that  $D_w(x_1) = \text{NW}_{\overline{L(M)}}(w, x_1)$ . The crucial observation that leads to the use of NW generator is that the desired pair  $(y'_w, z'_w)$  only depends on  $w|_{J_{x_1}}$ , which contains at most  $n$  bits of the input  $u \in \{0, 1\}^{n^{c/2}}$ . This is because  $w = a \cup u$  is a concatenation of a fixed  $a \in \{0, 1\}^{[n^c] \setminus J_{x_2}}$  and  $u$  viewed as  $u \in \{0, 1\}^{J_{x_2}}$ , which means that

$$w|_{J_{x_1}} = (a \cup u)|_{J_{x_1}} = a|_{J_{x_1}} \cup u|_{J_{x_1}},$$

where  $a|_{J_{x_1}}$  is fixed and  $u|_{J_{x_1}}$  only consists of the indices within  $J_{x_1} \cap J_{x_2}$  of size at most  $n$ .

As  $E(w)$  depends on at most  $n$  bits of the input  $u \in \{0, 1\}^{n^c}$ , we can implement it as a circuit that store all the answers for all  $2^n$  possibilities, which requires at most  $\text{poly}(2^n) = 2^{O(n)}$  gates. Concretely, the circuit works as follows: Given  $w \in \{0, 1\}^{n^c}$ , we firstly obtain  $u \in \{0, 1\}^{J_{x_2}}$  such that  $w = a \cup u$ ; let  $u' = u|_{J_{x_1}}$  be of length at most  $n$ , we look up the table to find the answer corresponding to  $u'$ .  $\square$

*Remark 5.12.* As in Remark 5.5, we note that the argument can be easily adapted to approximate any Boolean function  $g$  defined over  $n^k$  bits computable by a nondeterministic circuit of size  $2^{n^{k\delta}}$  using a deterministic circuit of size  $2^{O(n)}$ , instead of for just  $L(M) \cap \{0, 1\}^{n^k}$ . The provability of a circuit lower bound for a single language  $L(M)$  provides non-trivial circuits for any such  $g$ .

Based on this, we can also prove that under the same assumption (i.e., the provability of worst-case circuit lower bound in  $\text{TPV}$ ), for every constant  $\varepsilon \in (0, 1)$ ,  $s = s(m) = 2^{m^{o(1)}}$ , and sufficiently large  $m$ , any Boolean function  $g : \{0, 1\}^m \rightarrow \{0, 1\}$  that can be computable by a nondeterministic circuit of size  $s$  can also be approximated by a co-nondeterministic circuit  $D$  of size  $2^{m^\varepsilon}$ , that is:

$$\Pr_{x \sim \{0, 1\}^m} [C(x) = D(x)] \geq \frac{1}{2} + \frac{1}{2^{m^\varepsilon}}.$$

This can be done by setting  $k = \lceil 20/\varepsilon \rceil$ , padding dammy inputs to  $g : \{0, 1\}^m \rightarrow \{0, 1\}$  to obtain  $g' : \{0, 1\}^{m'} \rightarrow \{0, 1\}$ , where  $m' = \lceil m^{1/k} \rceil^k \leq 2m$  for sufficiently large  $m$ , and applying the observation above to  $g'$  with  $n = \lceil m^{1/k} \rceil$ .

### 5.2.3 Sketch of the general case

We now sketch how the argument presented in Section 5.2.2 can be generalised to the case that the Student-Teacher protocol runs for  $\ell \geq 3$  rounds. Recall that sets  $S_{x_1}, S_{x_1}^{\text{mist}}, S_{x_1, x_2}$  in Section 5.2.2 are defined as

$$\begin{aligned} S_{x_1} &\triangleq \{w \in \{0, 1\}^{n^c} \mid f_1(D_w) = x_1\} \\ S_{x_1}^{\text{mist}} &\triangleq \{w \in S_{x_1} \mid D_w(x_1) \neq M(x_1)\} \\ S_{x_1, x_2} &\triangleq \{w \in S_{x_1} \setminus S_{x_1}^{\text{mist}} \mid f_2(D_w, y'_w, z'_w) = x_2\} \end{aligned}$$

In the general case, we will define a sequence of  $x_1, x_2, \dots, x_\ell \in \{0, 1\}^n$  as well as the sets

$$S_1, S_1^{\text{mist}} \subseteq S_1, S_2 \subseteq S_1 \setminus S_1^{\text{mist}}, S_2^{\text{mist}} \subseteq S_2, \dots, S_\ell \subseteq S_{\ell-1} \setminus S_{\ell-1}^{\text{mist}}, S_\ell^{\text{mist}} \subseteq S_\ell.$$

For instance, if  $\ell = 3$ , we proceed as follows.

- (i) We initially argue as in Section 5.2.2 with  $\ell = 2$ . In Case 1 (i.e.,  $|S_{x_1}^{\text{mist}}| > (2/3) \cdot |S_{x_1}|$ ), we can simply apply the aforementioned circuit  $B_1$  to approximate  $L(M)$ . However, we can no longer conclude in its Case 2 (i.e.,  $|S_{x_1}^{\text{mist}}| < (2/3) \cdot |S_{x_1}|$ ) that  $x_2$  is a mistake of  $D_w$  for every  $w \in S_{x_1, x_2}$ . To address this, we define the set

$$S_{x_1, x_2}^{\text{mist}} \triangleq \{w \in S_{x_1, x_2} \mid D_w(x_2) \neq M(x_2)\} \subseteq S_{x_1, x_2}$$

and consider its density in  $S_{x_1, x_2}$ .

- (ii) If  $|S_{x_1, x_2}^{\text{mist}}|/|S_{x_1, x_2}| \geq 2/3$ , we know that for at least a  $2/3$  fraction of  $w \in S_{x_1, x_2}$ ,  $x_2$  is a mistake of  $D_w$ . As in Case 1 of Section 5.2.2, we apply Lemma E.1 (instead of a direct counting argument in Fact 5.11) to find a “good”  $a \in \{0, 1\}^{[n^c] \setminus J_{x_2}}$  such that  $S_{x_1, x_2} \upharpoonright_a / 2^{n^{c/2}} \geq \Omega(2^{-2n})$  and the density of  $S_{x_1, x_2}^{\text{mist}} \upharpoonright_a$  in  $S_{x_1, x_2} \upharpoonright_a$  is at least  $2/3 - 1/100$ . By plugging in this  $a$  into the circuit  $B_2$ , we will achieve agreement  $\geq 1/2 + \Omega(2^{-2n})$  with  $L(M)$ .
- (iii) Otherwise, we assume that  $|S_{x_1, x_2}^{\text{mist}}|/|S_{x_1, x_2}| < 2/3$ . Let  $(x_2, y_2, z_2) = f_2(D_w, y'_w, z'_w)$ . Similar to  $y'_w$  and  $z'_w$ , for every  $w \in S_{x_1, x_2} \setminus S_{x_1, x_2}^{\text{mist}}$ , we define  $(y''_w, z''_w)$  as the lexicographic first pair such that

$$\left[ M(x_2, y_2) = 0 \vee D_w(x_2, z_2) = 1 \right] \wedge \left[ M(x_2, y') = 1 \vee D_w(x_1, z') = 0 \right],$$

that is,  $(y''_w, z''_w)$  is the output of the canonical Teacher in the second round of the Student-Teacher protocol. Since  $S_{x_1, x_2}$  has density at least  $\Omega(2^{-2n})$ , we can find a string  $x_3 \in \{0, 1\}^n$  such that the following set

$$S_{x_1, x_2, x_3} \triangleq \{w \in S_{x_1, x_2} \setminus S_{x_1, x_2}^{\text{mist}} \mid f_3(D_w, y'_w, z'_w, y''_w, z''_w) = x_3\},$$

has density at least  $\Omega(2^{-3n})$ . Since  $x_3$  must be a mistake of  $D_w$  when  $\ell = 3$  and  $w \in S_{x_1, x_2, x_3}$ , and this set is sufficiently dense, we can obtain a deterministic circuit of size  $2^{O(n)}$  that achieves agreement  $\geq 1/2 + \Omega(2^{-3n})$  with  $L(M)$ .

The argument can be generalised in the natural way, which allows us to obtain a circuit of size  $2^{O(n)}$  that approximates  $L(M)$  with advantage  $\geq 1/2 + \Omega(2^{-\ell n})$  in the case of a disjunction of length  $\ell$  in the application of the KPT Witnessing (see Theorem 3.11). This deterministic circuit computes  $L(M)$  on inputs of length  $n^{c/2}$ , where  $c$  is an arbitrary constant.

*Remark 5.13.* Note that the approach breaks down in theories where the number of rounds in the Student-Teacher game obtained from Theorem 3.11 is polynomial in the relevant parameter, as in the case of Buss’s theory  $S_2^1$  (see, e.g., [Kra92]). In the latter case, one can get up to  $\ell = \text{poly}(2^n)$  rounds in the corresponding witnessing theorem, and the advantage of the resulting deterministic circuit under a naive extension of the presented proof becomes trivial.

### 5.3 Extensions of the technique and unprovability of weaker lower bounds

As noted in [PS21], one can use hardness amplification to weaken the average-case hardness in the unprovability result (Theorem 5.1). By an adaptation of the proof of Theorem 5.1 via Remarks 5.5 and 5.12 and an application of Theorem 3.5, we can obtain the following unprovability result.

**Theorem 5.14.** *For every  $n_0 \in \mathbb{N}$  and  $\delta \in \mathbb{Q} \cap (0, 1)$ , if  $M$  is a nondeterministic machine whose running time is bounded by some constructive function  $t(n) = 2^{n^{o(1)}}$ , then<sup>31</sup>*

$$\mathsf{TPV} \not\vdash \mathsf{LB}(M, s, m, n_0),$$

where  $s(n) = 2^{n^\delta}$  and  $m(n) = 2^n/n$ .

As a consequence, for every language  $L \in \mathsf{NTIME}[2^{n^{o(1)}}]$  and  $\delta > 0$  it is consistent with  $\mathsf{TPV}$  that there are infinitely many input lengths  $n$  and a co-nondeterministic circuit  $D_n$  of size  $\leq 2^{n^\delta}$  such that

$$\Pr_{x \sim \{0,1\}^n} [L(x) = D_n(x)] \geq 1 - 1/n.$$

*Proof of Theorem 5.14.* Let  $n_0, \delta, M, s(n) = 2^{n^\delta}$ , and  $m(n) = 2^n/n$  be as above. Assume towards a contradiction that

$$\mathsf{TPV} \vdash \mathsf{LB}(M, s, m, n_0).$$

Let  $L \triangleq L(M)$  be the language defined by  $M$ . We argue as follows.

- (i) Under the provability of an almost-everywhere average-case lower bound against conondeterministic circuits, it follows by the soundness of  $\mathsf{TPV}$  that (in the standard model) for every sequence  $\{E_n\}_{n \geq 1}$  of *deterministic* circuits  $E_n$  of size  $\leq 2^{n^\delta}$ , if  $n \geq n_0$  then

$$\Pr_{x \sim \{0,1\}^n} [L(x) = E_n(x)] \leq 1 - 1/n.$$

- (ii) From the provability of  $\mathsf{LB}(M, s, m, n_0)$ , it follows that  $\mathsf{TPV}$  proves the sentence  $\mathsf{LB}_{\text{wst}}(M, s, n_0)$  which states a *worst-case* lower bound for  $M$  against conondeterministic circuits of the same size. By adapting the argument presented in Section 5.2 (see Remarks 5.5 and 5.12), the provability of  $\mathsf{LB}_{\text{wst}}(M, s, n_0)$  in  $\mathsf{TPV}$  implies that, in the standard model, for *every* sequence  $\{g_n\}_{n \geq 1}$  of functions in  $\mathsf{NSIZE}[2^{n^{o(1)}}]$ ,  $\varepsilon > 0$ , and large enough  $n$ , there is a deterministic circuit  $C'$  defined over  $n$  input variables and of size  $2^{n^\varepsilon}$  such that

$$\Pr_{x \sim \{0,1\}^n} [g_n(x) = C'(x)] \geq 1/2 + 2^{-n^\varepsilon}. \quad (10)$$

- (iii) Let  $\{f_n\}_{n \geq 1}$  be the sequence of functions in  $\mathsf{NTIME}[2^{n^{o(1)}}]$  obtained from  $L$ , i.e.,  $f(x) = 1$  if and only if  $x \in L$ . Note that this sequence satisfies the hypothesis of Theorem 3.5 for  $s_1(n) = 2^{n^{o(1)}}$  and  $s_2(n) = 2^{n^\delta}$  for sufficiently large  $n$ . Let  $\{h_m\}_{m \geq 1}$  be the sequence of functions in  $\mathsf{NSIZE}[2^{m^{o(1)}}]$  obtained by an application of this result, we know that for sufficiently large  $n$  and any deterministic circuit  $C$  of size  $(2^{m^{\gamma\delta}})^\gamma$ , it holds that

$$\Pr_{x \sim \{0,1\}^m} [h_m(x) = C(x)] \leq 1/2 + 2^{-\gamma m^{\gamma\delta}}.$$

Now the hardness of  $h_m$  according to Theorem 3.5 contradicts the upper bound provided in Equation (10), if we take  $\varepsilon = (1/2) \cdot \delta \cdot \gamma$  and consider large enough input lengths.

This shows that  $\mathsf{TPV} \not\vdash \mathsf{LB}(M, s, m, n_0)$ , as desired. □

---

<sup>31</sup>The original statement in [PS21] is slightly weaker: they require the nondeterministic machine  $M$  to be in polynomial-time instead of  $t(n)$  time. We obtain such quantitative improvement by explicitly computing the complexity overhead of the hardness amplification in [HVV06] (see Theorem 3.5).



## 6 Unprovability of Strong Complexity Lower Bounds in Bounded Arithmetic

In this section, we establish the unprovability of strong  $\Sigma_i^p$ -vs- $\Pi_i^p$ -style lower bounds in bounded arithmetic. Our result generalises a previous unprovability result from [PS21] in two directions: (1) it holds for stronger theories  $T_{PV}^i$  instead of only  $T_{PV}^1$ ; and (2) the lower bound sentence in our unprovability result is more natural in the sense that the hard problem is quantified within the theory, instead of in the meta-theory.

Due to the complexity of the argument, we will first show in Section 6.1 how to generalise the unprovability result in [PS21] to  $T_{PV}^i$ . Then in Section 6.2 we combine this extension with the new game-theoretic witnessing theorem and with other ideas to obtain our main result, which has both features mentioned above.

### 6.1 Unprovability of lower bounds in expressive theories

For  $i \geq 1$ , recall that  $T_{PV}^i$  is the theory consisting of all true (in the standard model)  $\forall \Sigma_{i-1}^b(PV)$  sentences. For instance,  $T_{PV}^1$  is the universal true theory of PV. We want to generalize the unprovability of strong nondeterministic circuit lower bounds in  $T_{PV}^1$  to  $T_{PV}^i$  for all  $i \geq 1$ , stated as follows.<sup>32</sup>

**Theorem 6.1.** *Fix  $i \geq 1$ . Let  $t(n) = 2^{n^{o(1)}}$  be a constructive time bound, and  $M$  be a  $\Pi_i$ -TIME[ $t(n)$ ] machine and  $LB^i(M, s, m, n_0)$  be the  $\mathcal{L}_{PV}$ -sentence: for all  $n \in \text{LogLog}$  with  $n > n_0$  and  $C \in \Sigma_i\text{-SIZE}[s(n)]$ , there exist  $m$  distinct inputs  $x_1, \dots, x_m$  such that  $M(x_j) \neq C(x_j)$  for all  $j \in [m]$ . Then*

$$T_{PV}^i \not\vdash LB^i(M, s, m, n_0)$$

for  $s(n) = 2^{n^\delta}$ ,  $m(n) = 2^n/2 - 2^n/2^{n^\delta}$ , and  $\delta \in \mathbb{Q} \cap (0, 1)$ .

To obtain the unprovability of strong complexity lower bounds, we rely on a witnessing theorem that extracts computational information from a proof of the lower bound sentence  $LB^i(M, s, m, n_0)$ . We discuss the quantifier complexity of (the worst-case complexity analogue of) the  $LB^i(M, s, m, n_0)$  sentence in Section 6.1.1. As its formalization results in a  $\forall \Sigma_{i+1}^b(PV)$  sentence, note that when  $i > 1$  we can no longer directly apply the KPT Witnessing Theorem, as in Section 5. (In addition, for  $i > 1$  the theory  $T_{PV}^i$  is not universal, which is needed when applying this result.) A key aspect of our argument is to introduce an appropriate universal theory with the right abstractions and term complexity (see Section 3.5).

#### 6.1.1 Witnessing for $\Pi_i$ vs $\Sigma_i$ lower bounds

Let  $LB_{\text{wst}}^i(M, s, n_0)$  be the following *worst-case* lower bound sentence in the language  $\mathcal{L}_{PV}$ :

*For all  $n \in \text{LogLog}$  with  $n > n_0$  and circuit  $D \in \Sigma_i\text{-SIZE}[s(n)]$ , there exists an input  $x$  of length  $n$ , such that  $D(x) \neq M(x)$ .*

More formally, we have

$$LB_{\text{wst}}^i(M, s, n_0) \triangleq \forall n \in \text{LogLog with } n > n_0, \forall \text{ circuit } D \in \Sigma_i\text{-SIZE}[s(n)] \\ \exists x \in \{0, 1\}^n \text{ such that } \text{Error}(D, x),$$

<sup>32</sup>While in Section 5 we considered a lower bound for a nondeterministic machine against co-nondeterministic circuits, it will be more convenient for us in this section to phrase the statement as  $\Pi_i$ -machines against  $\Sigma_i$ -circuits. Note that this is inconsequential, as the results are equivalent via complementation.

where  $\text{Error}(D, x)$  is a sentence stating that  $M(x) \neq D(x)$ . Since  $M$  is a  $\Pi_i$ -machine and  $D$  is a  $\Sigma_i$ -circuit, in the language  $\mathcal{L}_{\text{PV}}$ , the sentence  $\phi_1(D, x) \triangleq (M(x) = 1 \wedge D(x) = 0)$  is in  $\Pi_i^b$  and the sentence  $\phi_2(D, x) \triangleq (M(x) = 0 \wedge D(x) = 1)$  is in  $\Sigma_i^b$ .

**Fact 6.2.** *Let  $m(n) \geq 1$ . If  $\text{T}_{\text{PV}}^i \vdash \text{LB}^i(M, s, m, n_0)$  then  $\text{T}_{\text{PV}}^i \vdash \text{LB}_{\text{wst}}^i(M, s, n_0)$ .*

*Proof.* This is immediate for any reasonable formalization of the sentence  $\text{LB}^i(M, s, m, n_0)$ , since it states an average-case lower bound (at least  $m(n) \geq 1$  mistakes) while  $\text{LB}_{\text{wst}}^i(M, s, n_0)$  states a worst-case lower bound (i.e. at least one mistake).  $\square$

Assume that

$$\begin{aligned}\phi_1(D, x) &\triangleq \forall y \in \{0, 1\}^{O(s(n))} \phi'_1(D, x, y), \\ \phi_2(D, x) &\triangleq \exists z \in \{0, 1\}^{O(s(n))} \phi'_2(D, x, z),\end{aligned}$$

for some  $\Sigma_{i-1}^b$ -formula  $\phi'_1$  and  $\Pi_{i-1}^b$ -formula  $\phi'_2$ , respectively. Note that the lengths of the strings  $y$  and  $z$  are bounded by  $O(s(n))$  since we obtain from them parts of the computation of the circuit  $D$  (of size  $s(n)$ ) and of the machine  $M$  (with running time  $2^{n^{o(1)}} < s(n)$ ). Then  $\text{Error}(D, x) \triangleq \phi_1(D, x) \vee \phi_2(D, x)$  is logically equivalent to the formula

$$\text{Error}'(D, x) \triangleq \exists z \in \{0, 1\}^{O(s(n))} \forall y \in \{0, 1\}^{O(s(n))} (\phi'_1(D, x, y) \vee \phi'_2(D, x, z)).$$

Next, consider the universal theories  $\text{U}_{\text{PV}}^i$  and  $\text{UT}_{\text{PV}}^i$  introduced in Section 3.5.

**Lemma 6.3.** *Let  $\text{ULB}_{\text{wst}}^i(M, s, n_0)$  be a  $\Pi_3^b$ -sentence in  $\mathcal{L}(\text{U}_{\text{PV}}^i)$  defined as follows:*

$$\begin{aligned}\text{ULB}_{\text{wst}}^i(M, s, n_0) &\triangleq \forall n \in \text{LogLog with } n \geq n_0, \forall \text{ circuit } D \in \Sigma_i\text{-SIZE}[s(n)] \\ &\quad \exists x \in \{0, 1\}^n \exists z \in \{0, 1\}^{O(s(n))} \\ &\quad \forall y \in \{0, 1\}^{O(s(n))} (f_{\phi'_1}(D, x, y) = 1 \vee f_{\phi'_2}(D, x, z) = 1).\end{aligned}$$

*Then  $\text{U}_{\text{PV}}^i$  proves  $\text{LB}_{\text{wst}}^i(M, s, n_0) \leftrightarrow \text{ULB}_{\text{wst}}^i(M, s, n_0)$ . Moreover,  $\text{UT}_{\text{PV}}^i$  proves  $\text{LB}_{\text{wst}}^i(M, s, n_0) \leftrightarrow \text{ULB}_{\text{wst}}^i(M, s, n_0)$ .*

*Proof.* By the discussion above,  $\text{LB}_{\text{wst}}^i(M, s, n_0)$  is logically equivalent to

$$\begin{aligned}&\forall n \in \text{LogLog with } n \geq n_0, \forall \text{ circuit } D \in \Sigma_i\text{-SIZE}[s(n)] \\ &\quad \exists x \in \{0, 1\}^n \exists z \in \{0, 1\}^{O(s(n))} \\ &\quad \forall y \in \{0, 1\}^{O(s(n))} (\phi'_1(D, x, y) = 1 \vee \phi'_2(D, x, z) = 1),\end{aligned}$$

which is further equivalent to  $\text{ULB}_{\text{wst}}^i(M, s, n_0)$  in  $\text{U}_{\text{PV}}^i$  by Lemma 3.17. The provability of the same sentence in  $\text{UT}_{\text{PV}}^i$  follows from Theorem 3.22.  $\square$

Note that the  $\mathcal{L}(\text{U}_{\text{PV}}^i)$ -sentence  $\text{ULB}_{\text{wst}}^i(M, s, n_0)$  has low quantifier complexity. By exploring the connection between  $\text{T}_{\text{PV}}^i$  and the universal theory  $\text{UT}_{\text{PV}}^i$ , we can show the following witnessing result.

**Lemma 6.4** (Witnessing lemma for  $\text{LB}(M, s, m, n_0)$ ). *Let  $i \geq 1$ ,  $M$  be a  $\Pi_i\text{-TIME}[2^{n^{o(1)}}]$  machine,  $\delta \in (0, 1)$ ,  $n_0 \geq 1$ ,  $s(n) = 2^{n^\delta}$ , and  $m(n) = 2^n/2 - 2^n/2^{n^\delta}$ . If  $\text{T}_{\text{PV}}^i \vdash \text{LB}^i(M, s, m, n_0)$ , then there exist  $\ell \in \mathbb{N}$  and  $\ell$  algorithms  $A_1, A_2, \dots, A_\ell$  such that:*

- Every  $A_i$  is computable in  $\text{FP}^{\Sigma_{i-1}^p}$  over inputs of length of order  $N = 2^n$ .
- For every  $i \in [\ell]$ , the input of  $A_i$  consists of  $1^N$ ,  $1^n$  for  $n = \log N$ , an  $n$ -input circuit  $D \in \Sigma_i\text{-SIZE}[s(n)]$ , and  $i - 1$  strings  $y_1, \dots, y_{i-1}$ ; the output of  $A_i$  is a pair  $(x_i, z_i) \in \{0, 1\}^n \times \{0, 1\}^{O(s(n))}$ .<sup>33</sup>
- Let  $h : (n, D, x) \mapsto y$  be the following function. Given  $n$ , a string  $x \in \{0, 1\}^n$ , and a circuit  $D \in \Sigma_i\text{-SIZE}[s(n)]$ , output a  $y$  such that  $\neg \phi'_1(D, x, y)$  if such  $y \in \{0, 1\}^{O(s(n))}$  exists, or 0 otherwise.
- For all  $n > n_0$  and circuit  $D \in \Sigma_i\text{-SIZE}[s(n)]$ , let

$$\begin{aligned} (x_1, z_1) &\triangleq A_1(1^n, D) & y_1 &\triangleq h(n, D, x_1) \\ (x_2, z_2) &\triangleq A_2(1^n, D, y_1) & y_2 &\triangleq h(n, D, x_2) \\ &\vdots & &\vdots \\ (x_\ell, z_\ell) &\triangleq A_\ell(1^n, D, y_1, \dots, y_{\ell-1}) & y_\ell &\triangleq h(n, D, x_\ell). \end{aligned}$$

Then there is an index  $v \in [\ell]$  such that  $D(x_v) \neq M(x_v)$ .

*Proof.* Let  $i, M, \delta, n_0, s(n), m(n)$  be defined as above. Assume that  $\text{T}_{\text{PV}}^i \vdash \text{LB}^i(M, s, m, n_0)$ . Then, by Fact 6.2, it follows that  $\text{T}_{\text{PV}}^i \vdash \text{LB}_{\text{wst}}^i(M, s, m, n_0)$ . Using Theorem 3.18 and Lemma 6.3, we get that  $\text{UT}_{\text{PV}}^i \vdash \text{ULB}_{\text{wst}}^i(M, s, m, n_0)$ .

Since  $\text{ULB}_{\text{wst}}^i(M, s, m, n_0)$  is a  $\forall \Sigma_2^b$ -sentence and  $\text{UT}_{\text{PV}}^i$  is a universal theory, we can invoke the KPT witnessing theorem (Theorem 3.11) to obtain constantly many  $\mathcal{L}_{\text{PV}}^i$ -terms  $A_1, A_2, \dots, A_\ell$  witnessing the existential quantifier given counter-examples to the innermost universal quantifier. By Theorem 3.21 and using that  $n \in \text{LogLog}$ , each  $A_i$  is computable in  $\text{FP}^{\Sigma_{i-1}^p}$  over an input of order  $N = 2^n$ . Furthermore, since the function  $h$  is a valid counter-example oracle for the innermost universal quantifier, it is easy to check that the conclusion of the lemma follows from the guarantee provided by KPT witnessing.  $\square$

### 6.1.2 Proof of Theorem 6.1

**Theorem** (Theorem 6.1, restated). *Fix  $i \geq 1$ . Let  $M$  be a  $\Pi_i\text{-TIME}[2^{n^{o(1)}}]$  machine and  $\text{LB}^i(M, s, m, n_0)$  be the  $\mathcal{L}_{\text{PV}}$ -sentence: for all  $n \in \text{LogLog}$  with  $n > n_0$  and  $\Sigma_i\text{-SIZE}[s(n)]$ -circuit  $C$ , there exist  $m$  distinct inputs  $x_1, \dots, x_m$  such that  $M(x_j) \neq C(x_j)$  for all  $j \in [m]$ . Then*

$$\text{T}_{\text{PV}}^i \not\vdash \text{LB}^i(M, s, m, n_0)$$

for  $s(n) = 2^{n^\delta}$ ,  $m(n) = 2^n/2 - 2^n/2^{n^\delta}$ , and  $\delta \in \mathbb{Q} \cap (0, 1)$ .

*Proof.* Suppose that  $\text{T}_{\text{PV}}^i \vdash \text{LB}^i(M, s, m, n_0)$  for some  $M \in \Pi_i\text{-TIME}[2^{n^{o(1)}}]$ ,  $n_0 \in \mathbb{N}$ ,  $s(n) = 2^{n^\delta}$ ,  $m(n) = 2^n/2 - 2^n/2^{n^\delta}$ , and  $\delta \in \mathbb{Q} \cap (0, 1)$ , there exist an  $\ell \in \mathbb{N}$  and  $\ell$  algorithms  $A_1, A_2, \dots, A_\ell$  as described by Lemma 6.4. Similar to [PS21], we will utilize the algorithms  $A_i$  to show that  $M$  can be non-trivially approximated by  $\Sigma_i\text{-SIZE}[2^{n^\delta}]$  circuits for some  $n > n_0$ , leading to a contradiction to the soundness of  $\text{T}_{\text{PV}}^i$ .

Let  $c$  be a constant to be determined later, and  $\text{NW}_f(w, x)$  be the Nisan-Wigderson generator with seed length  $|w| = n^c$ , output length  $2^n$ , “hard” function  $f : \{0, 1\}^{n^{c/2}} \rightarrow \{0, 1\}$  (therefore each subset in the combinatorial design has size  $n^{c/2}$ ),  $|x| = n$ , and any two distinct subsets in the combinatorial design with

<sup>33</sup>Formally, since  $n \in \text{LogLog}$  in our formalization, each  $A_i$  has access to an input  $\alpha$  of length  $|\alpha| = N = 2^n$ . For convenience of notation, when discussing  $A_1, \dots, A_\ell$  we often omit the input  $1^N$  and concentrate on  $n$ , which is the key parameter.

intersection of size at most  $n$ . For every seed  $w \in \{0, 1\}^{n^c}$ , let  $D_w : \{0, 1\}^n \rightarrow \{0, 1\}$  be a  $\Sigma_i$ -circuit computing  $D_w(x) \triangleq \text{NW}_{\overline{M}}(w, x)$ , which is of size at most  $2^{n^{o(c)}} \leq 2^{n^\delta}$  for sufficient large  $n$ . We will find some  $w \in \{0, 1\}^{n^c}$  and use  $D_w$  as  $C$  in Lemma 6.4 to obtain a  $\Sigma_i\text{-SIZE}[2^{O(n)}]$  circuit  $B$  approximating  $M$  on input length  $n^{c/2}$ , i.e.,  $\Pr_{u \in \{0, 1\}^{n^{c/2}}} [B(u) = M(u)] \geq \frac{1}{2} + 2^{-O(n)}$ . Then by choosing  $c > 2/\delta$  and sufficiently large  $n$ , we can prove the theorem.

**Case 1.** Recall that in Lemma 6.4,  $A_1$  takes  $1^n$  and an  $n$ -input circuit  $D \in \Sigma_i\text{-SIZE}[s(n)]$  as input and output a pair  $(x, y) \in \{0, 1\}^n \times \{0, 1\}^{O(s(n))}$ . By an averaging argument, there is an  $x_1 \in \{0, 1\}^n$  such that for a uniformly random  $w \in \{0, 1\}^{n^c}$ , with probability at least  $2^{-n}$ ,  $A_1(1^n, D_w)$  outputs  $(x_1, \cdot)$ . Fix this  $x_1$  and let

$$S_1 \triangleq \left\{ w \in \{0, 1\}^{n^c} \mid A_1(1^n, D_w) = (x_1, \cdot) \right\},$$

$$S_1^{\text{mist}} \triangleq \left\{ w \in S_1 \mid D_w(x_1) \neq M(x_1) \right\}.$$

By the definition of  $x_1$ , we know that  $|S_1|/2^{n^c} \geq 2^{-n}$ .

In Case 1 we assume that  $|S_1^{\text{mist}}| > (2/3) \cdot |S_1|$ , handling the other scenario in a subsequent case. For any  $w \in \{0, 1\}^{n^c}$ , we know that  $D_w(x_1) = \text{NW}_{\overline{M}}(w, x_1) = \overline{M}(w|_{J_{x_1}})$ , where  $J_{x_1}$  is the subset of indices corresponding to the  $x_1$ -th row of the combinatorial design. By Lemma E.1, there is an assignment  $a \in \{0, 1\}^{[n^c] \setminus J_{x_1}}$  for the indices outside of  $J_{x_1}$  such that  $|S_1 \upharpoonright_a|/2^{n^{c/2}} \geq 2^{-O(n)}$  and  $|S_1^{\text{mist}} \upharpoonright_a|/|S_1 \upharpoonright_a| \geq 3/5$ . Fix an  $a \in \{0, 1\}^{[n^c] \setminus J_{x_1}}$  as above. Let  $b_1 \triangleq M(x_1) \in \{0, 1\}$ . We define a randomized circuit  $B_1 : \{0, 1\}^{n^{c/2}} \times \{0, 1\} \rightarrow \{0, 1\}$ , where the second input is regarded as a random bit, as follows (see Algorithm 5 and recall the notation from Section 3.2).

**Input :** The input  $u \in \{0, 1\}^{n^{c/2}}$  for  $M$  and a bit  $r \in \{0, 1\}$

**Advice:**  $x_1 \in \{0, 1\}^n$ ,  $a \in \{0, 1\}^{[n^c] \setminus J_{x_1}}$ , and  $b_1 = M(x_1)$

- 1 Let  $w = r_{x_1}(a, u)$  and  $(x, \cdot) = A_1(1^n, D_w)$ ;
- 2 If  $x \neq x_1$ , **return**  $r$ ;
- 3 Otherwise, **return**  $b_1$ .

**Algorithm 5:** Randomized Circuit  $B_1$  for  $M$

We first analyse the complexity of  $B_1$ . Let  $m = n^{c/2} = |u|$  be the input length. Since  $A_1$  is computable in  $\text{FP}^{\Sigma_{i-1}^P}$ , it is easy to see that  $B_1 \in \text{SIZE}^{\Sigma_{i-1}^P}[2^{O(n)}]$ . By Theorem 3.2, we get that  $B_1 \in \Sigma_i\text{-SIZE}[2^{O(n)}]$ . So we only need to show that for an random bit  $r \in \{0, 1\}$ ,  $B(x, r)$  approximates  $M(x)$  well.

For any input  $u \in \{0, 1\}^{n^{c/2}}$  such that  $u \in S_1 \upharpoonright_a$ , we know that

$$\begin{aligned} B(u, r) = M(u) &\iff M(x_1) = M(u) && (x = x_1 \text{ by the definition of } S_1, B(u, r) = b_1 = M(x_1)) \\ &\iff M(x_1) \neq D_w(x_1) && (D_w(x_1) = \text{NW}_{\overline{M}}(w, x_1) = \overline{M}(w|_{J_{x_1}}) = \overline{M}(u)) \\ &\iff u \in S_1^{\text{mist}} \upharpoonright_a. \end{aligned}$$

This means that  $B(u, r)$  and  $M$  agree on at least  $3/5$  of the inputs  $u \in S_1 \upharpoonright_a$ . In the other case, the circuit  $B$  outputs the random bit  $r$ , therefore for some fixed bit  $r^* \in \{0, 1\}$ ,  $B_1(u, r^*)$  and  $M(u)$  agree on at least  $1/2$  of the inputs  $u \notin S_1 \upharpoonright_a$ . Since  $|S_1 \upharpoonright_a|/2^{n^{c/2}} \geq 2^{-O(n)}$ , we obtain that

$$\Pr_{u \in \{0, 1\}^{n^{c/2}}} [B_1(u, r^*) = M(u)] \geq \frac{3}{5} \cdot \frac{|S_1 \upharpoonright_a|}{2^{n^{c/2}}} + \frac{1}{2} \cdot \left(1 - \frac{|S_1 \upharpoonright_a|}{2^{n^{c/2}}}\right) = \frac{1}{2} + 2^{-O(n)}.$$

**Case 2.** Assume that  $|S_1^{\text{mist}}| \leq (2/3) \cdot |S_1|$  instead. For every  $w \in S_1$ , we define  $y_1(w) \triangleq h(n, D_w, x_1)$  to be the output of the counter-example oracle  $h$  in Lemma 6.4. Again by an averaging argument, there is an  $x_2 \in \{0, 1\}^n$  such that for a uniformly random  $w \in S_1 \setminus S_1^{\text{mist}}$ , with probability at least  $2^{-n}$ ,  $A_2(1^n, D_w, y_1(w)) = (x_2, \cdot)$ . Fix this  $x_2$ . Let  $S_2$  and  $S_2^{\text{mist}}$  be the sets defined as follows:

$$S_2 \triangleq \left\{ w \in S_1 \setminus S_1^{\text{mist}} \mid A_2(1^n, D_w, y_1(w)) = (x_2, \cdot) \right\},$$

$$S_2^{\text{mist}} \triangleq \{ w \in S_2 \mid D_w(x_2) \neq M(x_2) \}.$$

By the definition of  $x_2$  we know that  $|S_2|/2^{n^c} \geq (1/3) \cdot 2^{-O(n)} = 2^{-O(n)}$ .

In this case, we further assume that  $|S_2^{\text{mist}}| > (2/3) \cdot |S_2|$ . By construction, for any  $w \in \{0, 1\}^{n^c}$ ,  $D_w(x_2) = \overline{M}(w|_{J_{x_2}})$ . By Lemma E.1, there is an assignment  $a \in \{0, 1\}^{[n^c] \setminus J_{x_2}}$  for the indices outside of  $J_{x_2}$  such that  $|S_2 \upharpoonright_a|/2^{n^{c/2}} \geq 2^{-O(n)}$  and  $|S_2^{\text{mist}} \upharpoonright_a|/|S_2 \upharpoonright_a| \geq 3/5$ . Fix this  $a$ . We will need the following subroutine to complete this case.

( $\nabla$ ) Given  $w \in S_1$  of the form  $a \cup u$  ( $u \in \{0, 1\}^{J_{x_2}}$ ), there is a deterministic circuit  $E(w)$  of size at most  $2^{O(n)}$  that outputs  $(y_1(w), e_1(w))$ , where  $y_1(w) = h(n, D_w, x_1)$  and  $e_1(w) \in \{0, 1\}$  such that  $e_1(w) = 1$  if and only if  $w \in S_1^{\text{mist}}$ .

Note that the circuit  $E(w)$  is used to simulate the counter-example oracle  $h$  in the first round of the KPT-style game. Let  $b_2 \triangleq M(x_2)$ . We construct a randomized circuit  $B_2$  as follows (see Algorithm 6), discussing the claim ( $\nabla$ ) later in the proof.

**Input :** The input  $u \in \{0, 1\}^{n^{c/2}}$  for  $M$  and random bit  $r \in \{0, 1\}$   
**Advice:**  $x_1, x_2 \in \{0, 1\}^n$ ,  $a \in \{0, 1\}^{[n^c] \setminus J_{x_2}}$  as discussed,  $b_2 = M(x_2)$ , and  $\Gamma$  to support the subroutine ( $\nabla$ )

- 1 Let  $w = r_{x_2}(a, u)$  and  $(\hat{x}_1, \hat{z}_1) = A_1(1^n, D_w)$ ;
- 2 If  $\hat{x}_1 \neq x_1$ , then **return the random bit**  $r$ ; // after this step,  $w \in S_1$
- 3 Let  $(y_1(w), e_1(w)) = E(w)$  by ( $\nabla$ );
- 4 If  $e_1(w) = 1$ , then **return the random bit**  $r$ ; // after this step,  $w \in S_1 \setminus S_1^{\text{mist}}$
- 5 Let  $(\hat{x}_2, \hat{z}_2) = A_2(1^n, D_w, y_1(w))$ ;
- 6 If  $\hat{x}_2 \neq x_2$ , then **return the random bit**  $r$ ;
- 7 Otherwise, **return**  $b_2$ . // reaching this line if and only if  $w \in S_2$

**Algorithm 6:** Randomized circuit  $B_2$  for  $M$

Let  $m = n^{c/2} = |u|$  be the input length. We first show that  $B_2 : \{0, 1\}^m \times \{0, 1\} \rightarrow \{0, 1\}$  can be implemented by a  $\Sigma_i$ -circuit with size  $2^{O(n)}$ . Both  $A_1$  and  $A_2$  are  $\text{FP}^{\Sigma_{i-1}^p}$  algorithms with input length  $\text{poly}(2^n)$ , so both of them can be implemented by  $\text{SIZE}^{\Sigma_{i-1}^p}[2^{O(n)}]$  circuits. By ( $\nabla$ ) we also know that  $E(w)$  can be implemented by a  $2^{O(n)}$ -size circuit. As a result,  $B_2 \in \text{SIZE}^{\Sigma_{i-1}^p}[2^{O(n)}] \subseteq \Sigma_i\text{-SIZE}[2^{O(n)}]$ , where the last inclusion follows from Theorem 3.2.

Now we prove the correctness of the algorithm  $B_2$ . By construction, it is easy to verify that the algorithm reaches the last line if and only if  $w = r_{x_2}(a, u) \in S_2$ . Therefore  $B_2$  will output a random bit when  $w \notin S_2$  (i.e.,  $u \notin S_2 \upharpoonright_a$ ) and output  $b_2$  when  $w \in S_2$  (i.e.,  $u \in S_2 \upharpoonright_a$ ). In the former case,  $B_2$  agrees with  $M$  on  $1/2$  of the inputs for an  $r^* \in \{0, 1\}$ , which will be hard-wired into the circuit. In the latter case, with probability at least  $3/5$  over  $u \in \{0, 1\}^{n^{c/2}}$ ,  $u \in S_2^{\text{mist}} \upharpoonright_a$ , which further means that

$$M(u) = M(w|_{J_{x_2}}) = \overline{D}_w(x_2) = M(x_2) = b_2 = B_2(u, r).$$

Since  $|S_2 \upharpoonright_a|/2^{n^{c/2}} \geq 2^{-O(n)}$ , we can conclude that  $B_2(u, r^*)$  agrees with  $M(u)$  on  $\frac{1}{2} + 2^{-O(n)}$  of the inputs  $u \in \{0, 1\}^{n^{c/2}}$ .

**Case  $j \geq 2$ .** Using the technique for Case 2, we can in fact deal with all the remaining cases. Let  $j \in \{2, 3, \dots, \ell\}$ . We define the following notations recursively:

- (i)  $y_{j-1}(w) \triangleq h(n, D_w, x_{j-1})$ .
- (ii)  $x_j \in \{0, 1\}^n$  be the lexicographically first string such that for a uniformly random string  $w \in S_{j-1} \setminus S_{j-1}^{\text{mist}}$ , with probability at least  $2^{-n}$ ,  $A_j(1^n, D_w, y_1(w), \dots, y_{j-1}(w)) = (x_j, \cdot)$ .
- (iii) Define  $S_j$  and  $S_j^{\text{mist}}$  as the sets

$$S_j \triangleq \left\{ w \in S_{j-1} \setminus S_{j-1}^{\text{mist}} \mid A_j(1^n, D_w, y_1(w), \dots, y_{j-1}(w)) = (x_j, \cdot) \right\}$$

$$S_j^{\text{mist}} \triangleq \{w \in S_j \mid D_w(x_j) \neq M(x_j)\}.$$

In Case  $j \geq 2$  we assume that (1)  $|S_j^{\text{mist}}| > (2/3) \cdot |S_j|$ , and (2) for every  $i \in \{1, 2, \dots, j-1\}$ ,  $|S_i^{\text{mist}}|/|S_i| \leq 2/3$ . Crucially, by the definition of each of these sets and the conclusion of Lemma 6.4, we get that by reaching  $j = \ell$  we necessarily have  $S_\ell = S_\ell^{\text{mist}}$ , so the case analysis is complete.

The following lemma will be needed later in the proof.

**Lemma 6.5.** *For every  $1 \leq i < j$ , we have  $x_i \neq x_j$ .*

*Proof.* Suppose that  $x_i = x_j$  for  $i < j$ . First, it follows from the construction that  $S_j \cap S_i^{\text{mist}} = \emptyset$ . Therefore, for any  $w \in S_j^{\text{mist}} \subseteq S_j$ , we have  $M(x_i) = D_w(x_i)$ . On the other hand, by definition, for any  $w \in S_j^{\text{mist}}$ , we have  $M(x_j) \neq D_w(x_j)$ . Note that the assumption of Case  $j \geq 2$  implies that  $S_j^{\text{mist}}$  is nonempty. Take any  $w^* \in S_j^{\text{mist}}$ . Under the hypothesis that  $x_i = x_j$ , the previous claims yield that both  $M(x_i) = D_{w^*}(x_i)$  and  $M(x_i) \neq D_{w^*}(x_i)$ , which is contradictory.  $\square$

Under assumptions (1) and (2), one can prove by induction that  $|S_j|/2^{n^c} = 2^{-O(n)}$ , therefore by Lemma E.1, there is an assignment  $a \in \{0, 1\}^{[n^c] \setminus J_{x_j}}$  such that  $|S_j \upharpoonright_a|/2^{n^{c/2}} \geq 2^{-O(n)}$  and  $|S_j^{\text{mist}} \upharpoonright_a|/|S_j \upharpoonright_a| \geq 3/5$ . Similarly to  $(\nabla)$  in Case 2, we need the following computation  $(\nabla_j^i)$  for every  $i \in \{1, 2, \dots, j-1\}$ .

$(\nabla_j^i)$  Given  $w \in S_i$  of the form  $a \cup u$  ( $u \in \{0, 1\}^{J_{x_j}}$ ), there is a deterministic circuit  $E_i(w)$  of size at most  $2^{O(n)}$  that outputs  $(y_i(w), e_i(w))$ , where  $y_i(w) = h(n, D_w, x_i)$  and  $e_i(w) \in \{0, 1\}$  such that  $e_i(w) = 1$  if and only if  $w \in S_i^{\text{mist}}$ .

Note that  $(\nabla_2^1)$  is simply  $(\nabla)$  in Case 2. With these subroutines we can construct a randomized circuit  $B_j$  that approximates  $M$  well as follows (see Algorithm 7).

Now we analyze the complexity and correctness of the algorithm  $B_j$ .

**(Complexity).** Let  $m = n^{c/2}$  be the input length. Since every  $A_i$  is computable in  $\text{FP}^{\Sigma_{i-1}^p}$  with input length  $\text{poly}(2^n, s(n)) = 2^{O(n)}$ , and every  $E_i$  is computable by a  $2^{O(n)}$ -size deterministic circuit, we know that  $B_j \in \text{SIZE}^{\Sigma_{j-1}^p}[2^{O(n)}] \subseteq \Sigma_j\text{-SIZE}[2^{O(n)}]$  (see Theorem 3.2).

**(Correctness).** The correctness of  $B_j$  is proved similarly to Case 2. As noted in the comments appearing in the pseudo-code, for any  $w = r_{x_j}(a, u)$  with  $u \in \{0, 1\}^{n^{c/2}}$ , we can prove by induction that the algorithm reaches the end of the  $i$ -th iteration within the for-loop if and only if  $w \in S_i \setminus S_i^{\text{mist}}$  for every

**Input :** The input  $u \in \{0, 1\}^{n^{c/2}}$  for  $M$  and random bit  $r \in \{0, 1\}$   
**Advice:**  $x_1, \dots, x_j \in \{0, 1\}^n$ ,  $a \in \{0, 1\}^{[n^c] \setminus J_{x_j}}$  as discussed above,  $b_j = M(x_j)$ , and  $\Gamma$  to support the subroutines  $(\nabla_j^i)$

```

1 Let  $w = r_{x_j}(a, u)$ ;
2 for  $i = 1, 2, \dots, j$  do
3   Let  $(\hat{x}_i, \hat{z}_i) = A_i(1^n, D_w, y_1, \dots, y_{i-1})$ ;
4   If  $\hat{x}_i \neq x_i$ , then return the random bit  $r$ ;
   // reaching this line iff  $w \in S_i$ 
5   if  $i < j$  then
6     Let  $(y_i(w), e_i(w)) = E_i(w)$  by  $(\nabla_j^i)$ ;
7     If  $e_i(w) = 1$ , then return the random bit  $r$ ;
     // otherwise,  $x \in S_i \setminus S_i^{\text{mist}}$ 
8   end
9 end
   // reaching this line iff  $w \in S_j$ 
10 return  $b_j$ ;

```

**Algorithm 7:** Randomized circuit  $B_j$  for  $M$

$i \in \{1, 2, \dots, j-1\}$ . Furthermore, on such inputs the algorithm reaches the last line if and only if  $w \in S_j$ . This means that, for an appropriate fixed bit  $r^* \in \{0, 1\}$ , on inputs of this form the algorithm agrees with  $M$  on at least  $1/2$  of  $w \notin S_j$  and on at least  $3/5$  of  $w \in S_j$ . Since  $|S_j \upharpoonright_a|/2^{n^{c/2}} = 2^{-O(n)}$ ,  $B_j(\cdot, r^*)$  achieves an advantage of  $2^{-O(n)}$  with  $M(\cdot)$ .

**Implementation of  $(\nabla)$ .** To finish the proof, we need to upper bound the circuit complexity of the computation  $E_i(w)$  in  $\nabla_j^i$  for  $1 \leq i < j \leq \ell$ , which is used to simulate the counter-example oracle  $h$  and to check if  $w \in S_i^{\text{mist}}$ , for  $w$  of the appropriate form. Let  $j \in \{2, 3, \dots, \ell\}$  and  $i \in \{1, 2, \dots, j-1\}$ . Thanks to Lemma 6.5, we know that  $x_i \neq x_j$ . Recall that  $D_w(x) \triangleq \text{NW}_{\overline{M}}(w, x)$ , where  $w \triangleq r_{x_j}(a, u)$ . Since any two distinct subsets in the combinatorial design of the Nisan-Wigderson generator have intersection size at most  $n$ , we get that  $|J_{x_i} \cap J_{x_j}| \leq n$ . Notice that for  $u \in \{0, 1\}^{n^{c/2}}$ ,  $h(n, D_w, x_i)$  for  $w = r_{x_j}(a, u)$  only depends on  $w|_{J_{x_i}}$ , which contains at most  $n$  bits of  $u$ . As a result, we can hard-wire the answers of all  $2^n$  cases and construct a  $2^{O(n)}$  circuit to compute  $y_i(w) = h(n, D_w, x_i)$ . The value  $e_i(w)$  can be computed in a similar way. Overall, we obtain that  $E_i(w) \triangleq (y_i(w), e_i(w))$  can be computed on all relevant inputs by a (non-uniform) deterministic circuit of size  $2^{O(n)}$ .

**Wrapping things up.** By the case analysis above, for every  $c \geq 2$  and every sufficiently large  $n$ , there always exists a  $\Sigma_i\text{-SIZE}[2^{O(n)}]$  circuit  $B$  with input length  $m = n^{c/2}$  such that

$$\Pr_{u \in \{0, 1\}^m} [B(u) = M(u)] \geq \frac{1}{2} + 2^{-O(n)}.$$

By taking  $c \geq 2/\delta$  we get, that for infinitely many values of  $n$ ,  $L(M) \cap \{0, 1\}^n$  can be approximated with advantage at least  $2^{-n^\delta}$  by  $\Sigma_i\text{-SIZE}[2^{n^\delta}]$  circuits. This leads to a contradiction, since under  $\text{T}_{\text{PV}}^i \vdash \text{LB}^i(M, s, m, n_0)$  and from the soundness of  $\text{T}_{\text{PV}}^i$  we obtain that, for sufficiently large  $n$ ,  $M$  cannot be approximated with advantage  $2^{-n^\delta}$  by  $\Sigma_i\text{-SIZE}[2^{n^\delta}]$  circuits.  $\square$



As pointed out in Remark 5.5 and Remark 5.12, we note that assuming the provability of *worst-case* circuit lower bound, the approximation of the machine  $M \in \Pi_i\text{-TIME}[2^{n^{o(1)}}]$  by deterministic  $\Sigma_{i-1}^p$ -oracle circuits of small size also works for any sequence  $\{f_n\}_{n \geq 1}$  of functions computable in  $\Pi_i\text{-SIZE}[2^{n^{o(1)}}]$ . Concretely, instead of defining  $D_w(x) \triangleq \text{NW}_{\overline{M}}(w, x)$ , we define  $D_w(x) \triangleq \text{NW}_{f'}(w, x)$  for  $f'(x) \triangleq \neg f(x)$ , and proceed the argument as above. By padding dammy input bits as in Remark 5.12, we obtain the following corollary.

**Corollary 6.6.** *Fix  $i \geq 1$ . Let  $M$  be a  $\Pi_i\text{-TIME}[2^{n^{o(1)}}]$  machine and  $\text{LB}_{\text{wst}}^i(M, s, n_0)$  be the worst-case lower bound sentence defined as above. Assume that for some  $\delta \in (0, 1) \cap \mathbb{Q}$  and  $s(n) = 2^{n^\delta}$ ,  $\text{T}_{\text{PV}}^i$  proves  $\text{LB}_w^i(M, s, n_0)$ . Then for every constant  $\varepsilon \in (0, 1)$ , every sufficiently large  $n$ , and circuit  $C \in \Pi_i\text{-SIZE}[2^{n^\delta}]$ , there is a  $\Sigma_{i-1}^p$ -oracle circuit  $D$  of size  $2^{n^\varepsilon}$  such that*

$$\Pr_{x \sim \{0,1\}^n} [C(x) = D(x)] \geq \frac{1}{2} + \frac{1}{2^{n^\varepsilon}}.$$

### 6.1.3 Relaxing the average-case complexity parameter

Recall that in Section 5.3, we showed how to obtain the unprovability of circuit lower bounds with a weaker average-case complexity parameter via hardness amplification. This will also be the case here, since the hardness amplification in [HVV06] generalises to all levels in the polynomial hierarchy (see Theorem 3.5).

**Theorem 6.7.** *Fix  $i \geq 1$ . Let  $M$  be a  $\Pi_i\text{-TIME}[t(n)]$  machine for some constructive function  $t(n) = 2^{n^{o(1)}}$  and  $\text{LB}^i(M, s, m, n_0)$  be defined as in Theorem 6.1. Then for every constant  $\delta \in \mathbb{Q} \cap (0, 1)$ ,  $s(n) \triangleq 2^{n^\delta}$ ,  $m(n) \triangleq 2^n/n$ , and  $n_0 \in \mathbb{N}$ ,  $\text{T}_{\text{PV}}^i \not\vdash \text{LB}^i(M, s, m, n_0)$ .*

*Proof.* Towards a contradiction, we assume that  $\text{T}_{\text{PV}}^i \vdash \text{LB}^i(M, s, m, n_0)$  and argue as follows.

- (i) Under the provability of the almost-everywhere average-case lower bound  $\text{LB}(M, s, m, n_0)$ , we obtain from the soundness of  $\text{T}_{\text{PV}}^i$  that (in the standard model) for every sequence  $\{E_n\}_{n \geq 1}$  of  $\Sigma_{i-1}^p$ -oracle circuits of size  $\leq 2^{n^\delta}$  and  $n \geq n_0$ , we have

$$\Pr_{x \sim \{0,1\}^n} [M(x) = E_n(x)] \leq 1 - \frac{1}{n}.$$

- (ii) From the provability of  $\text{LB}(M, s, m, n_0)$ , under reasonable formalization, we can also show that  $\text{LB}_{\text{wst}}^i(M, s, n_0)$  is provable in  $\text{T}_{\text{PV}}^i$ . We then get from Corollary 6.6 that for every  $\varepsilon \in (0, 1)$ , every sufficiently large  $n$ , and circuit  $C \in \Pi_i\text{-SIZE}[2^{n^\delta}]$ , there is a  $\Sigma_{i-1}^p$ -oracle circuit  $D$  of size  $2^{n^\varepsilon}$  such that

$$\Pr_{x \sim \{0,1\}^n} [C(x) = D(x)] \geq \frac{1}{2} + \frac{1}{2^{n^\varepsilon}}. \quad (11)$$

- (iii) Assume that  $n$  is sufficiently large and  $f_n : \{0, 1\}^n \rightarrow \{0, 1\}$  is defined as  $f(x) = M(x)$ . Note that this function satisfies the hypothesis of Theorem 3.5 for  $s_1(n) = 2^{n^{o(1)}}$  and  $s_2(n) = 2^{n^\delta}$ , hence we can obtain a function  $h_\ell : \{0, 1\}^\ell \rightarrow \{0, 1\}$  for some  $\ell = O(n^2)$  such that for every  $\Sigma_{i-1}^p$ -oracle circuit  $D$  of size  $2^{\gamma \ell^{\gamma \delta}}$ , it holds that

$$\Pr_{x \in \{0,1\}^\ell} [h_\ell(x) = D(x)] \leq \frac{1}{2} + \frac{1}{2^{\gamma \ell^{\gamma \delta}}}.$$

By setting  $\varepsilon = (1/2) \cdot \delta \cdot \gamma$ , this violates the upper bound in Equation (11) when  $n$  is sufficiently large.

As a result, we know that  $\text{TPV}^i \not\models \text{LB}^i(M, s, m, n_0)$  for every  $i \geq 1$ .  $\square$

## 6.2 Unprovability of lower bound sentences of higher quantifier complexity

In this section, we extend the unprovability results to sentences of higher quantifier complexity that formalize separations between non-uniform circuit classes. Recall that  $\Sigma_i\text{-SIZE}[s(n)]$  and  $\Pi_i\text{-SIZE}[s(n)]$  refer to  $\Sigma_i$ -circuits and  $\Pi_i$ -circuits of size  $s(n)$ , respectively. Let  $\text{LB}^i(s_1, s_2, m, n_0)$  denote the following  $\mathcal{L}_{\text{PV}}$ -sentence:

$$\begin{aligned} & \forall n \in \text{LogLog with } n \geq n_0 \exists C \in \Pi_i\text{-SIZE}[s_1(n)] \forall D \in \Sigma_i\text{-SIZE}[s_2(n)] \\ & \exists m = m(n) \text{ distinct } n\text{-bit strings } x^1, \dots, x^m \text{ s.t. Error}(C, D, x^i) \text{ for all } i \in [m], \end{aligned}$$

where  $\text{Error}(C, D, x)$  means that the circuits  $C$  and  $D$  do not agree on the input  $x$ . It's easy to see that  $\text{Error}(C, D, x)$  is a disjunction of a  $\Sigma_i^b$ -formula and a  $\Pi_i^b$ -formula. Observe that, already for  $i = 1$ ,  $\text{LB}^i(s_1, s_2, m, n_0)$  is a  $\forall\Sigma_4^b$ -sentence.

**Theorem 6.8.** *For every  $i \geq 1$ ,  $n_0 \in \mathbb{N}$ ,  $\delta \in \mathbb{Q} \cap (0, 1)$  and  $d \geq 1$ ,  $\text{TPV}^i \not\models \text{LB}^i(s_1, s_2, m, n_0)$ , where  $s_1(n) = n^d$ ,  $s_2(n) = 2^{n^\delta}$  and  $m(n) = 2^n/2 - 2^n/2^{n^\delta}$ .*

### 6.2.1 Witnessing lemma for lower bound sentences

Similar to the technique we used in the previous section, we need to apply the witnessing theorem to the lower bound sentences. We define the worst-case version of this lower bound to be the following formula  $\text{LB}_{\text{wst}}^i(s_1, s_2, n_0)$ .

$$\begin{aligned} & \forall n \in \text{LogLog with } n \geq n_0 \exists C \in \Pi_i\text{-SIZE}[s_1(n)] \forall D \in \Sigma_i\text{-SIZE}[s_2(n)] \\ & \exists x \in \{0, 1\}^n \text{ s.t. Error}(C, D, x). \end{aligned}$$

Let  $\phi_1(C, D, x) \triangleq (C(x) = 1 \wedge D(x) = 0)$  be a  $\Pi_i^b$ -formula and  $\phi_2(C, D, x) \triangleq (C(x) = 0 \wedge D(x) = 1)$  be a  $\Sigma_i^b$ -formula. Note that  $\text{Error}(C, D, x) \triangleq \phi_1(C, D, x) \vee \phi_2(C, D, x)$ . Assume that

$$\begin{aligned} \phi_1(C, D, x) & \triangleq \forall y \in \{0, 1\}^{O(s(n))} \phi'_1(C, D, x, y), \\ \phi_2(C, D, x) & \triangleq \exists z \in \{0, 1\}^{O(s(n))} \phi'_2(C, D, x, z), \end{aligned}$$

where  $\phi'_1$  is a  $\Sigma_{i-1}^b$ -formula and  $\phi'_2$  is a  $\Pi_{i-1}^b$ -formula. Note that the lengths of  $y$  and  $z$  are bounded by  $O(s(n))$  since they are parts of the computation of the circuits  $C$  and  $D$ .

**Lemma 6.9.** *Let  $\text{ULB}_{\text{wst}}^i(s_1, s_2, n_0)$  be a  $\forall\Sigma_4^b$ -sentence in  $\mathcal{L}(\text{UPV}^i)$  defined as follows:*

$$\begin{aligned} \text{ULB}_{\text{wst}}^i(s_1, s_2, n_0) & \triangleq \forall n \in \text{LogLog with } n \geq n_0, \exists \text{ circuit } C \in \Pi_i\text{-SIZE}[s_1(n)] \\ & \forall \text{ circuit } D \in \Sigma_i\text{-SIZE}[s_2(n)], \\ & \exists x \in \{0, 1\}^n \exists z \in \{0, 1\}^{O(s(n))} \forall y \in \{0, 1\}^{O(s(n))} \\ & \left( f_{\phi'_1}(C, D, x, y) = 1 \vee f_{\phi'_2}(C, D, x, z) = 1 \right). \end{aligned}$$

*Then  $\text{UPV}^i$  proves  $\text{LB}_{\text{wst}}^i(s_1, s_2, n_0) \leftrightarrow \text{ULB}_{\text{wst}}^i(s_1, s_2, n_0)$ . Moreover,  $\text{UT}_{\text{PV}}^i$  proves  $\text{LB}_{\text{wst}}^i(s_1, s_2, n_0) \leftrightarrow \text{ULB}_{\text{wst}}^i(s_1, s_2, n_0)$ .*

*Proof.* The provability in  $\text{U}_{\text{PV}}^i$  follows from the provability of the defining axioms for  $f_\alpha$  (see Lemma 3.17). In turn, the provability in  $\text{UT}_{\text{PV}}^i$  follows from Theorem 3.22.  $\square$

**Lemma 6.10.** *Assume that  $\text{T}_{\text{PV}}^i \vdash \text{LB}^i(s_1, s_2, m, n_0)$ . There is an integer  $\ell \in \mathbb{N}$  and  $\text{FP}^{\Sigma_{i-1}^p}$  algorithms  $P_1, Q_1, P_2, Q_2, \dots, P_\ell, Q_\ell$  such that the following condition holds.<sup>34</sup>*

*Let  $n > n_0$ ,  $g$  be a function that maps a  $\Pi_i\text{-SIZE}[s_1(n)]$ -circuit to a  $\Sigma_i\text{-SIZE}[s_2(n)]$  circuit, and  $D_C \triangleq g(C)$ . Let  $h : (n, C, D, x) \mapsto y$  be the function such that  $y$  is the lexicographic first string in  $\{0, 1\}^{O(s(n))}$  such that  $\neg \phi_1^t(C, D, x, y)$  holds or 0 if such a string does not exist. Let*

$$\begin{array}{lll} P_1(1^n) = C_1 & Q_1(1^n, D_{C_1}) = (x_1, z_1) & h(n, C_1, D_{C_1}, x_1) = y_1 \\ P_2(1^n, D_{C_1}, y_1) = C_2 & Q_2(1^n, D_{C_1}, D_{C_2}, y_1) = (x_2, z_2) & h(n, C_2, D_{C_2}, x_2) = y_2 \\ \vdots & \vdots & \vdots \\ P_\ell(1^n, D_{C_{1\dots\ell-1}}, y_{1\dots\ell-1}) = C_\ell & Q_\ell(1^n, D_{C_{1\dots\ell-1}}, y_{1\dots\ell-1}) = (x_\ell, z_\ell) & h(n, C_\ell, D_{C_\ell}, x_\ell) = y_\ell. \end{array}$$

*Then there is  $k \in [\ell]$  such that  $\text{Error}(C_k, D_{C_k}, x_k)$  holds.*

*Proof.* Suppose that  $\text{T}_{\text{PV}}^i \vdash \text{LB}^i(s_1, s_2, m, n_0)$ . Then we also have  $\text{T}_{\text{PV}}^i \vdash \text{LB}_{\text{wst}}^i(s_1, s_2, n_0)$ , which further means by Theorem 3.18, Lemma 6.9, and Theorem 3.22 that  $\text{UT}_{\text{PV}}^i \vdash \text{ULB}_{\text{wst}}^i(s_1, s_2, n_0)$ . Recall that  $\text{UT}_{\text{PV}}^i$  is a universal theory closed under if-then-else (see Theorem 3.22). By Theorem 4.20, there are an  $\ell \in \mathbb{N}$  and a sequence of  $\ell$   $\mathcal{L}$ -strategies  $\tau_1^t, \tau_2^t, \dots, \tau_\ell^t$  for the truthifier such that for any fixed strategy  $\tau^f$  of the falsifier, at least one of the strategies beats  $\tau^f$  in  $\ell$  sequential plays of the evaluation game with  $\tau_1^t, \tau_2^t, \dots, \tau_\ell^t$  vs  $\tau^f$ .

In particular, consider the following strategy for the falsifier: if the truthifier chooses  $C$  in the first round of the game, the falsifier will choose  $D_C$ ; then if the truthifier chooses  $x, z$  in the second round, the falsifier will choose  $h(n, C, D_C, x)$ . It is easy to see that the claim in the lemma is precisely the winning property of  $\tau_1^t, \dots, \tau_\ell^t$  against this particular strategy for the falsifier, given the corresponding auxiliary information.  $\square$

### 6.2.2 Proof of Theorem 6.8

**Theorem (Reminder of Theorem 6.8).** *For every  $i \geq 1$ ,  $n_0 \in \mathbb{N}$ ,  $\delta \in \mathbb{Q} \cap (0, 1)$  and  $d \geq 1$ ,  $\text{T}_{\text{PV}}^i \not\vdash \text{LB}^i(s_1, s_2, m, n_0)$ , where  $s_1(n) = n^d$ ,  $s_2(n) = 2^{n^\delta}$  and  $m(n) = 2^n/2 - 2^n/2^{n^\delta}$ .*

*Proof.* Assume that  $\text{T}_{\text{PV}}^i \vdash \text{LB}^i(s_1, s_2, m, n_0)$ . We will derive a contradiction to the soundness of  $\text{T}_{\text{PV}}^i$  by showing that for sufficiently large  $n$  and all  $\Pi_i$ -circuits  $M : \{0, 1\}^{n^{c/2}} \rightarrow \{0, 1\}$  of size  $s_1(n^{c/2})$ , there is a  $\Sigma_i$ -circuit  $B : \{0, 1\}^{n^{c/2}} \rightarrow \{0, 1\}$  of size at most  $s_2(n^{c/2})$  that agrees with  $M$  on all but at most  $m(n^{c/2})$  inputs, for some constant  $c \in \mathbb{N}$  which will be determined later.

Let  $\text{NW}_f(w, x)$  be the Nisan-Wigderson generator with:  $f : \{0, 1\}^{n^{c/2}} \rightarrow \{0, 1\}$ , seed length  $|w| = n^c$ ,  $|x| = n + n^d$ , and any two distinct subsets in the combinatorial design of intersection of size at most  $O(n^d)$ . Designs with these parameters are known to exist (see Section 3.2).

By Lemma 6.10, we have  $\ell \in \mathbb{N}$  and  $\text{FP}^{\Sigma_{i-1}^p}$  machines  $P_1, P_2, \dots, P_\ell, Q_1, \dots, Q_\ell$  as described. Let  $M : \{0, 1\}^{n^{c/2}} \rightarrow \{0, 1\}$  be a  $\Pi_i$ -circuit of size  $s_1(n^{c/2})$  as described above. Let  $D_{w,C} : \{0, 1\}^n \rightarrow \{0, 1\}$  be a  $\Sigma_i$ -circuit of size at most  $s_2(n)$  computing  $D_{w,C}(x) \triangleq \text{NW}_{\overline{M}}(w, x \parallel C)$  for  $w \in \{0, 1\}^{n^c}$  and  $C \in \{0, 1\}^{n^d}$ .<sup>35</sup> We would like to find some suitable  $w$  and apply Lemma 6.10 with  $g : C \mapsto D_{w,C}$  to obtain a circuit  $B \in \text{SIZE}^{\Sigma_{i-1}^p}[2^{O(n^d)}] \subseteq \Sigma_i\text{-SIZE}[2^{O(n^d)}]$  approximating  $M$ , i.e.  $\Pr_u[B(u) = M(u)] \geq \frac{1}{2} + 2^{-O(n^d)}$ . By choosing  $c$  as a constant much larger than  $d$ , we can prove the theorem.

<sup>34</sup>As in the statement of Lemma 6.4, the input length of these algorithms is of order  $N = 2^n$ , since in our formalisation  $n \in \text{LogLog}$ . In order to be succinct, we simply write  $1^n$  as one of the inputs, since  $n$  is the key parameter for us.

<sup>35</sup>We use  $u \parallel v$  to denote the concatenation of binary strings  $u$  and  $v$ . Jumping ahead, the idea of concatenating  $x \parallel C$  when defining the NW generator will allow us to establish an analogue of Lemma 6.5 in this proof.

**Case 1.** Let  $C_1 \triangleq P_1(1^n)$ . By an averaging argument, there is an  $x_1 \in \{0, 1\}^n$  such that for a uniformly random  $w \in \{0, 1\}^{n^c}$ , with probability at least  $2^{-n}$ , the first coordinate of  $Q_1(1^n, D_{w, C_1})$  is  $x_1$ . Fix this  $x_1$  and let

$$S_1 \triangleq \left\{ w \in \{0, 1\}^{n^c} \mid Q_1(1^n, D_{w, C_1}) = (x_1, \cdot) \right\},$$

$$S_1^{\text{mist}} \triangleq \left\{ w \in S_1 \mid D_{w, C_1}(x_1) \neq C_1(x_1) \right\}.$$

By the definition of  $x_1$  we get that  $|S_1|/2^{n^c} \geq 2^{-n}$ .

In this case we assume that  $|S_1^{\text{mist}}| \geq (2/3) \cdot |S_1|$ , dealing with the other situation in a subsequent case analysis. For any  $w$ , we know that  $D_{w, C_1}(x_1) = \text{NW}_{\overline{M}}(w, x_1 \| C_1) = \overline{M}(w|_{J_{x_1 \| C_1}})$ , where  $J_{x_1 \| C_1}$  is the subset of indices corresponding to the  $(x_1 \| C_1)$ -th row of the combinatorial design. By Lemma E.1, there is an assignment  $a \in \{0, 1\}^{[n^c] \setminus J_{x_1 \| C_1}}$  for the indices outside of  $J_{x_1 \| C_1}$  such that  $|S_1 \upharpoonright_a|/2^{n^c/2} \geq 2^{-O(n)}$  and  $|S_1^{\text{mist}} \upharpoonright_a|/|S_1 \upharpoonright_a| \geq 3/5$ .

Now we fix  $a \in \{0, 1\}^{[n^c] \setminus J_{x_1 \| C_1}}$  as above. Let  $b_1 \triangleq C_1(x_1) \in \{0, 1\}$ . We define a randomized circuit  $B_1$  with access to a random bit  $r \in \{0, 1\}$  as follow (see Algorithm 8).

**Input :** The input  $u \in \{0, 1\}^{n^c/2}$  for  $M$  and random bit  $r \in \{0, 1\}$

**Advice:**  $x_1 \in \{0, 1\}^n$ ,  $C_1 = P_1(1^n)$ ,  $a \in \{0, 1\}^{[n^c] \setminus J_{x_1 \| C_1}}$  as discussed, and  $b_1 = C_1(x_1)$

- 1 Let  $w = r_{x_1 \| C_1}(a, u)$  and  $(x, \cdot) = Q_1(1^n, D_{w, C_1})$ ;
- 2 If  $x \neq x_1$ , **return** the random bit  $r$ ;
- 3 Otherwise, **return**  $b_1$ .

**Algorithm 8:** Randomized circuit  $B_1$  for  $M$

Since  $Q_1 \in \text{FP}^{\Sigma_{i-1}^p}$  and  $|D_{w, C_1}| = 2^{n^{o(1)}}$ , it is clear that  $B_1 \in \text{SIZE}^{\Sigma_{i-1}^p}[2^{O(n)}] \subseteq \Sigma_i\text{-SIZE}[2^{O(n^d)}]$ , so we only need to verify that the randomized circuit  $B_1$  approximates  $M$ . For  $u \in \{0, 1\}^{n^c/2}$  such that  $u \in S_1 \upharpoonright_a$ , we have that

$$\begin{aligned} B_1(u, r) = M(u) &\iff C_1(x_1) = M(u) \quad (x = x_1 \text{ by the definition of } S_1, B(u, r) = b_1 = C(x_1)) \\ &\iff C_1(x_1) \neq D_{w, C_1}(x_1) \quad (D_{w, C_1}(x_1) = \text{NW}_{\overline{M}}(w, x_1 \| C_1) = M(u)) \\ &\iff u \in S_1^{\text{mist}} \upharpoonright_a. \end{aligned}$$

Therefore  $B_1$  and  $M$  agree on at least  $3/5$  of the inputs  $u \in S_1 \upharpoonright_a$ . In the other case, the circuit  $B$  simply outputs the random bit  $r$ , therefore for a specific  $r^* \in \{0, 1\}$ ,  $B_1(u, r^*)$  and  $M(u)$  agree on at least  $1/2$  of the inputs  $u \notin S_1 \upharpoonright_a$ . Since  $|S_1 \upharpoonright_a|/2^{n^c/2} \geq 2^{-O(n)}$ , we obtain that

$$\Pr_{u \in \{0, 1\}^{n^c/2}} [B_1(u, r^*) = M(u)] \geq \frac{3}{5} \cdot \frac{|S_1 \upharpoonright_a|}{2^{n^c/2}} + \frac{1}{2} \cdot \left(1 - \frac{|S_1 \upharpoonright_a|}{2^{n^c/2}}\right) = \frac{1}{2} + 2^{-O(n)}.$$

**Case 2.** Assume that  $|S_1^{\text{mist}}| \leq (2/3) \cdot |S_1|$ . Let  $h(n, C, D, x_1)$  be the function described in Lemma 6.10. Let  $y_1(w) \triangleq h(n, C_1, D_{w, C_1}, x_1)$  and  $C_2^w = P_2(1^n, D_{w, C_1}, y_1(w))$ . Again, by an averaging argument, there are  $C_2 \in \{0, 1\}^{n^d}$  and  $x_2 \in \{0, 1\}^n$  such that for a uniformly random  $w \in S_1 \setminus S_1^{\text{mist}}$ , with probability at least  $2^{-O(n^d)}$ ,  $C_2 = C_2^w$  and  $Q_2(1^n, D_{w, C_1}, D_{w, C_2}, y_1(w)) = (x_2, \cdot)$ . Fix this  $C_2$  and  $x_2$ . Let  $S_2$  and  $S_2^{\text{mist}}$  be sets defined as follows:

$$S_2 \triangleq \left\{ w \in S_1 \setminus S_1^{\text{mist}} \mid C_2 = C_2^w \wedge Q_2(1^n, D_{w, C_1}, D_{w, C_2}, y_1(w)) = (x_2, \cdot) \right\}$$

$$S_2^{\text{mist}} \triangleq \left\{ w \in S_2 \mid D_{w, C_2}(x_2) \neq C_2(x_2) \right\}$$

By the definitions of  $C_2$  and  $x_2$ , we know that  $|S_2|/2^{n^c} \geq (1/3) \cdot 2^{-O(n^d)} = 2^{-O(n^d)}$ .

In this case we assume that  $|S_2^{\text{mist}}| \geq (2/3) \cdot |S_2|$ . Similarly to Case 1, for any  $w \in \{0, 1\}^{n^c}$ ,  $D_{w, C_2}(x_2) = \overline{M}(w|_{J_{x_2 \| C_2}})$ . By Lemma E.1, there is an assignment  $a \in \{0, 1\}^{[n^c] \setminus J_{x_2 \| C_2}}$  for the indices outside of  $J_{x_2 \| C_2}$  such that  $|S_2 \upharpoonright_a|/2^{n^{c/2}} \geq 2^{-O(n^d)}$  and  $|S_2^{\text{mist}} \upharpoonright_a|/|S_2 \upharpoonright_a| \geq 3/5$ . Fix this string  $a$ . We will assume the following computation is possible in order to complete this case, returning to it later on:

( $\nabla$ ) Given  $w \in S_1$  of the form  $a \cup u$  ( $u \in \{0, 1\}^{J_{x_2 \| C_2}}$ ), there is a deterministic circuit  $E(w)$  of size at most  $2^{O(n^d)}$  that outputs  $(y_1(w), e_1(w))$ , where  $y_1(w) = h(n, C_1, D_{w, C_1}, x_1)$  and  $e_1(w) \in \{0, 1\}$  such that  $e_1(w) = 1$  if and only if  $w \in S_1^{\text{mist}}$ .

Note that if  $w \in S_1 \setminus S_1^{\text{mist}}$ ,  $y_1(w)$  given by  $E(w)$  witnesses that  $\neg \text{Error}(C_1, D_{w, C_1}, x_1)$ . Let  $b_2 \triangleq C_2(x_2)$ . We construct a randomized circuit  $B_2$  as follows (see Algorithm 9).

**Input :** The input  $u \in \{0, 1\}^{n^{c/2}}$  for  $M$  and random bit  $r \in \{0, 1\}$   
**Advice:**  $x_1, x_2 \in \{0, 1\}^n$ ,  $C_1, C_2 \in \{0, 1\}^{n^d}$ ,  $a \in \{0, 1\}^{[n^c] \setminus J_{x_2 \| C_2}}$  as discussed,  $b_2 = C_2(x_2)$ , and  $\Gamma$  to support the subroutine ( $\nabla$ )

- 1 Let  $w = r_{x_2 \| C_2}(a, u)$  and  $(\hat{x}_1, \hat{z}_1) = Q_1(1^n, D_{w, C_1})$ ;
- 2 If  $\hat{x}_1 \neq x_1$ , then **return** the random bit  $r$ ; // after this step,  $w \in S_1$
- 3 Let  $(y_1(w), e_1(w)) = E(w)$  by ( $\nabla$ );
- 4 If  $e_1(w) = 1$ , then **return** the random bit  $r$ ; // after this step,  $w \in S_1 \setminus S_1^{\text{mist}}$
- 5 Let  $\hat{C}_2 = P_2(1^n, D_{w, C_1}, y_1(w))$  and  $(\hat{x}_2, \hat{z}_2) = Q_2(1^n, D_{w, C_1}, D_{w, C_2}, y_1(w))$ ;
- 6 If  $\hat{x}_2 \neq x_2$  or  $\hat{C}_2 \neq C_2$ , then **return** the random bit  $r$ ;
- 7 Otherwise, **return**  $b_2$ . // reaching this line if and only if  $w \in S_2$

**Algorithm 9:** Randomized circuit  $B_2$  for  $M$

First, we analyze the complexity of  $B_2$ . Since  $Q_1, P_2, Q_2 \in \text{FP}^{\Sigma_{i-1}^p}$  and the input length for each of them is of order  $2^{O(n)}$ , they can be implemented by circuits of size  $2^{O(n)}$  with  $\Sigma_{i-1}^p$  oracles. We need  $2^{O(n^d)}$  gates to support the computation ( $\nabla$ ). Therefore,  $B_2 \in \text{SIZE}^{\Sigma_{i-1}^p}[2^{O(n^d)}] \subseteq \Sigma_i\text{-SIZE}[2^{O(n^d)}]$ .

By construction, it is easy to verify that the algorithm reaches the last line if and only if  $w = r_{x_2 \| C_2}(a, u) \in S_2$ . Therefore  $B_2$  will output a random bit when  $w \notin S_2$  and output  $b_2$  when  $w \in S_2$ . In the former case,  $B_2$  agrees with  $M$  on  $1/2$  of the inputs for an  $r^* \in \{0, 1\}$ . In the latter case, with probability at least  $3/5$ ,  $w|_{J_{x_2 \| C_2}} \in S_2^{\text{mist}} \upharpoonright_a$ , which further means that

$$M(u) = M(w|_{J_{x_2 \| C_2}}) = \overline{D_{w, C_2}}(x_2) = C_2(x_2) = b_2 = B_2(u, r).$$

As a result, it follows that

$$\Pr_{u \in \{0, 1\}^{n^{c/2}}} [B_2(u, r^*) = M(u)] \geq \frac{3}{5} \cdot \frac{|S_2 \upharpoonright_a|}{2^{n^{c/2}}} + \frac{1}{2} \cdot \left(1 - \frac{|S_2 \upharpoonright_a|}{2^{n^{c/2}}}\right) = \frac{1}{2} + 2^{-O(n^d)}.$$

**Case  $j \geq 2$ .** Using the technique for Case 2, we can in fact deal with all the remaining cases. Let  $j \in \{2, 3, \dots, \ell\}$ . We recursively define the following values:

- (i)  $y_{j-1}(w) \triangleq h(n, C_{j-1}, D_{w, C_{j-1}}, x_{j-1})$ .
- (ii)  $C_j^w \triangleq P_j(1^n, D_{w, C_1}, \dots, D_{w, C_{j-1}}, y_1(w), \dots, y_{j-1}(w))$ .

- (iii) Let  $C_j \in \{0, 1\}^{n^d}$  be the lexicographical first string (encoding an circuit) such that for a uniformly random string  $w \in S_{j-1} \setminus S_{j-1}^{\text{mist}}$ , with probability at least  $2^{-O(n^d)}$ ,  $C_j^w = C_j$ . The existence of  $C_j$  follows from a counting argument.
- (iv) Let  $x_j \in \{0, 1\}^n$  be the lexicographical first string such that for a uniformly random string  $w \in (S_{j-1} \setminus S_{j-1}^{\text{mist}}) \cap \{w \in \{0, 1\}^{n^c} \mid C_j^w = C_j\}$ , with probability at least  $2^{-n}$ ,

$$Q_j(1^n, D_{w, C_1}, \dots, D_{w, C_j}, y_1(w), \dots, y_{j-1}(w)) = (x_j, \cdot).$$

Thus, for a uniformly random string  $w \in S_{j-1} \setminus S_{j-1}^{\text{mist}}$ , with probability at least  $2^{-O(n^d)} \cdot 2^{-n} = 2^{-O(n^d)}$ ,  $C_j^w = C_j$  and  $Q_j(1^n, D_{w, C_1}, \dots, D_{w, C_j}, y_1(w), \dots, y_{j-1}(w)) = (x_j, \cdot)$ .

- (v)  $S_j$  and  $S_j^{\text{mist}}$  be sets recursively defined as

$$\begin{aligned} S_j &\triangleq \left\{ w \in S_{j-1} \setminus S_{j-1}^{\text{mist}} \mid C_j^w = C_j \wedge \right. \\ &\quad \left. Q_j(1^n, D_{w, C_1}, \dots, D_{w, C_j}, y_1(w), \dots, y_{j-1}(w)) = (x_j, \cdot) \right\} \\ S_j^{\text{mist}} &\triangleq \{w \in S_j \mid D_{w, C_j}(x_j) \neq C_j(x_j)\} \end{aligned}$$

In Case  $j$  we will assume that (1)  $|S_j^{\text{mist}}|/|S_j| \geq 2/3$  and (2) for any  $i \in \{1, 2, \dots, j-1\}$ ,  $|S_i^{\text{mist}}|/|S_i| < 2/3$ . In particular, by Lemma 6.10 we know that if we reach  $j = \ell$  then  $S_\ell = S_\ell^{\text{mist}}$ , so all the cases can be resolved in this way.

The following lemma will be useful later in the proof.

**Lemma 6.11.** *For every  $1 \leq i < j$ , we have  $(C_i, x_i) \neq (C_j, x_j)$ .*

*Proof.* First, note that  $S_j \cap S_i^{\text{mist}} = \emptyset$ . Also, since we are in case  $j$ ,  $S_j^{\text{mist}} \neq \emptyset$ , given that  $|S_j^{\text{mist}}| \geq 2/3 \cdot |S_j|$  and the (inductively established) density lower bound for  $|S_j|$ . Now take any  $w^* \in S_j^{\text{mist}}$ , i.e.,

$$C_j(x_j) \neq D_{w^*, C_j}(x_j). \quad (12)$$

Since  $S_j^{\text{mist}} \subseteq S_j$  and  $S_j \cap S_i^{\text{mist}} = \emptyset$ , we have that  $w^* \notin S_i^{\text{mist}}$ , i.e.,

$$C_i(x_i) = D_{w^*, C_i}(x_i). \quad (13)$$

Now if we had  $(C_i, x_i) = (C_j, x_j)$ , this would be in contradiction with Equation (12) and Equation (13). Consequently, either  $C_i \neq C_j$  or  $x_i \neq x_j$ .  $\square$

We can prove by induction that  $|S_j|/2^{n^c} = 2^{-O(n^d)}$ , therefore by Lemma E.1, there is an assignment  $a \in \{0, 1\}^{[n^c] \setminus J_{x_j} \parallel C_j}$  such that  $|S_j \upharpoonright_a|/2^{n^c/2} \geq 2^{-O(n^d)}$  and  $|S_j^{\text{mist}} \upharpoonright_a|/|S_j \upharpoonright_a| \geq 3/5$ . Fix this string  $a$ . Similar to  $(\nabla)$  in Case 2, we need the following computation  $(\nabla_j^i)$  for every  $i \in \{1, 2, \dots, j-1\}$ .

$(\nabla_j^i)$  Given  $w \in S_i$  of the form  $a \cup u$  ( $u \in \{0, 1\}^{J_{x_j} \parallel C_j}$ ), there is a deterministic circuit  $E_i(w)$  of size at most  $2^{O(n^d)}$  that outputs  $(y_i(w), e_i(w))$ , where  $y_i(w) = h(n, C_i, D_{w, C_i}, x_i)$  and  $e_i(w) \in \{0, 1\}$  such that  $e_i(w) = 1$  if and only if  $w \in S_i^{\text{mist}}$ .

Note that  $(\nabla_2^1) = (\nabla)$  by definition. Let  $b_j \triangleq C_j(x_j)$ . Using the subroutines described above, We can now present a randomized circuit  $B_j$  that approximates  $M$  (see Algorithm 10).

**Input :** The input  $u \in \{0, 1\}^{n^{c/2}}$  for  $M$  and random bit  $r \in \{0, 1\}$   
**Advice:**  $x_1, \dots, x_j \in \{0, 1\}^n$ ,  $C_1, \dots, C_j \in \{0, 1\}^{n^d}$ ,  $a \in \{0, 1\}^{[n^c] \setminus J_{x_j} \parallel C_j}$  as discussed,  
 $b_j = C_j(x_j)$ , and  $\Gamma$  to support the subroutines  $(\nabla_j^i)$

```

1 Let  $w = r_{x_j}(a, u)$ ;
2 for  $i = 1, 2, \dots, j$  do
3   Let  $\hat{C}_i = P_i(1^n, D_{w, C_1}, \dots, D_{w, C_{i-1}}, y_1(w), \dots, y_{i-1}(w))$ ;
4   If  $\hat{C}_i \neq C_i$ , then return the random bit  $r$ ;
5   Let  $(\hat{x}_i, \hat{z}_i) = Q_i(1^n, D_{w, C_1}, \dots, D_{w, C_i}, y_1(w), \dots, y_{i-1}(w))$ ;
6   If  $\hat{x}_i \neq x_i$ , then return the random bit  $r$ ;
   // reaching this line iff  $w \in S_i$ 
7   if  $i < j$  then
8     Let  $(y_i(w), e_i(w)) = E_i(w)$  by  $(\nabla_j^i)$ ;
9     If  $e_i(w) = 1$ , then return the random bit  $r$ ;
     // otherwise,  $x \in S_i \setminus S_i^{\text{mist}}$ 
10  end
11 end
   // reaching this line iff  $w \in S_j$ 
12 return  $b_j$ ;

```

**Algorithm 10:** Randomized circuit  $B_j$  for  $M$

Similarly to Case 2, we can see that  $B_j \in \text{SIZE}^{\Sigma_{i-1}^p}[2^{O(n^d)}]$ . Now we prove the correctness of  $B_j$ . By the definition of  $S_i$ , we can prove by induction that the algorithm reaches the end of the  $i$ -th iteration within the for-loop if and only if  $w \in S_i \setminus S_i^{\text{mist}}$  for any  $i \in \{1, 2, \dots, j-1\}$ . We can further check that the algorithm reaches the last line if and only if  $w \in S_j$ . This means that, by fixing an appropriate bit  $r^*$  as the random bit, the algorithm agrees with  $M$  on at least  $1/2$  of  $w \notin S_j$  of the form  $w = r_{x_j}(a, u)$ , and on at least  $3/5$  of  $w \in S_j$  of the form  $w = r_{x_j}(a, u)$ . As before, this translates into a correlation of  $2^{-O(n^d)}$  over a random input  $u \in \{0, 1\}^{n^{c/2}}$  using the lower bound on the density of  $S_j \upharpoonright_a$ .

**Implementation of  $(\nabla)$ .** To complete the proof it is sufficient to show that  $(\nabla_j^i)$  in the  $j$ -th step is computable by  $2^{O(n^d)}$ -size circuits, for all  $j \in \{2, 3, \dots, \ell\}$  and  $1 \leq i < j$ . Fix any  $j \in \{2, 3, \dots, \ell\}$  and  $i < j$ . Recall that  $h(n, C_i, D_{w, C_i}, x_i)$  finds the minimal  $y_i$  such that  $\neg \phi'_1(C_i, D_{w, C_i}, x_i, y_i)$  holds if  $\neg \phi_1(C_i, D_{w, C_i}, x_i)$ , where  $\neg \phi_1(C_i, D_{w, C_i}, x_i)$  means that  $C_i(x_i) = 0 \vee D_{w, C_i}(x_i) = 1$ . In case  $C_i(x_i) = 0$ , we only need to hard-wire a witness of it, since  $C_i$  and  $x_i$  are fixed with respect to  $w$ .

Now we assume that  $C_i(x_i) = 1$ . By the definition of  $D_{w, C_i}$ , we know that  $D_{w, C_i}(x_i) = \text{NW}_{\overline{M}}(w, x_i \parallel C_i)$ , where  $w = a \cup u$  for  $a \in \{0, 1\}^{[n^c] \setminus J_{x_j} \parallel C_j}$ ,  $u \in \{0, 1\}^{J_{x_j} \parallel C_j}$ . By Lemma 6.11,  $(x_i, C_i) \neq (x_j, C_j)$ . Therefore, by the definition of the NW generator, for  $w = a \cup u$  with the input  $u \in \{0, 1\}^{J_{x_j} \parallel C_j}$ , the output  $D_{w, C_i}(x_i)$ , as well as the desired witness of the outer-most quantified variable in case that  $D_{w, C_i}(x_i) = 1$ , depends on at most  $O(n^d)$  bits of  $u$ . In such case, we can hard-wire all the  $2^{O(n^d)}$  answers with a deterministic  $2^{O(n^d)}$ -size circuit.

Similarly, it is not hard to show that the computation  $e_i(w)$  can also be implemented by a deterministic circuit of at most this size.



**Wrapping things up.** Finally we can combine all these facts to conclude this theorem. By assuming  $T_{PV}^i \vdash LB^i(s_1, s_2, m, n_0)$ , we proved that for sufficiently large  $n$  and all  $\Pi_i$ -circuits  $M : \{0, 1\}^{n^{c/2}} \rightarrow \{0, 1\}$  of size  $s_1(n^{c/2}) = n^{dc/2}$ , there is a deterministic circuit  $B : \{0, 1\}^{n^{c/2}} \rightarrow \{0, 1\}$  with  $\Sigma_{i-1}^p$  oracle gates of size at most  $2^{O(n^d)}$  that agrees with  $M$  on a  $1/2 + 2^{-O(n^d)}$  fraction of inputs. By Theorem 3.2, we know that  $B$  can be implemented by  $\Sigma_i$ -circuits of size  $2^{O(n^d)}$ . If we choose  $c > 2d/\delta$ ,  $B$  is of size  $\leq 2^{(n^{c/2})^\delta}$  and agrees with  $M$  on a  $\geq 1/2 + 2^{-(n^{c/2})^\delta}$  fraction of the inputs  $u \in \{0, 1\}^{n^{c/2}}$ , which means that  $\mathbb{N} \models \neg LB(s_1, s_2, m, n_0)$  for the corresponding choice of  $m(n)$ . This is a contradiction to the soundness of  $T_{PV}^i$ .  $\square$

Similarly to what was noted in Remark 5.12 and Corollary 6.6, the proof presented above shows that one can approximate every  $\Pi_i$ -SIZE $[2^{n^{o(1)}}]$  circuit  $M$  by small-size  $\Sigma_{i-1}^p$ -oracle circuits, assuming the provability of the worst-case circuit lower bound sentence  $LB_{wst}(s_1, s_2, n_0)$ . We simply use  $D_{w,C}(x) \triangleq NW_{\overline{M}}(w, x \| C)$  and proceed as above. By padding dummy input bits as in Remark 5.12, we can obtain the following corollary.

**Corollary 6.12.** Fix  $i \geq 1$ . Assume that for some  $n_0 \in \mathbb{N}$ ,  $\delta \in \mathbb{Q} \cap (0, 1)$ , and  $d \geq 1$ ,  $T_{PV}^i \vdash LB_{wst}^i(s_1, s_2, n_0)$  for  $s_1(n) = n^d$  and  $s_2(n) = 2^{n^\delta}$ . Then for every constant  $\varepsilon > 0$ , every sufficiently large  $n$ , and circuit  $A \in \Pi_i$ -SIZE $[t(n)]$  where  $t(n) = 2^{n^{o(1)}}$  is some constructive function, there is a  $\Sigma_{i-1}^p$ -oracle circuit  $B$  of size  $2^{n^\varepsilon}$  such that

$$\Pr_{x \sim \{0,1\}^n} [A(x) = B(x)] \geq \frac{1}{2} + \frac{1}{2^{n^\varepsilon}}.$$

### 6.2.3 Relaxing the average-case complexity parameter

As in Section 6.1.3, we now utilize the hardness amplification theorem (see Theorem 3.5) to relax the average-case complexity parameter.

**Theorem 6.13.** For every  $i \geq 1$ ,  $n_0 \in \mathbb{N}$ ,  $\delta \in \mathbb{Q} \cap (0, 1)$ , and  $d \geq 1$ ,  $T_{PV}^i \not\vdash LB^i(s_1, s_2, m, n_0)$ , where  $s_1(n) = n^d$ ,  $s_2(n) = 2^{n^\delta}$ , and  $m = 2^n/n$ .

*Proof.* Suppose that  $s_1 = s_1(n)$ ,  $s_2 = s_2(n)$ ,  $m$ , and  $n_0$  are defined as above. Towards a contradiction, we assume that  $T_{PV}^i \vdash LB^i(s_1, s_2, m, n_0)$ .

- (i) Under the unprovability of the almost-everywhere average-case lower bound  $LB(s_1, s_2, m_0)$ , we obtain from the soundness of  $T_{PV}^i$  that (in the standard model) for sufficiently large  $n$ , there is a circuit  $C \in \Pi_i$ -SIZE $[s_1(n)]$  such that for every  $\Sigma_{i-1}^p$ -oracle circuit  $D$  of size  $2^{n^\delta}$ , we have

$$\Pr_{x \sim \{0,1\}^n} [C(x) = D(x)] \leq 1 - \frac{1}{n}.$$

- (ii) By the assumption that  $T_{PV}^i \vdash LB^i(s_1, s_2, m, n_0)$ , under any reasonable formalization, we know that  $T_{PV}^i$  also proves the *worst-case* version of the lower bound  $LB_{wst}^i(s_1, s_2, n_0)$ . Then by Corollary 6.12, we get that for every constant  $\varepsilon > 0$ , every sufficiently large  $n$ , and every  $\Pi_i$ -SIZE $[2^{n^{o(1)}}]$  circuit, there is a  $\Sigma_{i-1}^p$ -oracle circuit  $D$  of size  $2^{n^\varepsilon}$  such that

$$\Pr_{x \sim \{0,1\}^n} [C(x) = D(x)] \geq \frac{1}{2} + \frac{1}{2^{n^\varepsilon}}. \tag{14}$$

(iii) Now we assume that  $n$  is sufficiently large and  $f_n : \{0, 1\}^n \rightarrow \{0, 1\}$  is the function computable by  $\Pi_i\text{-SIZE}[s_1(n)]$  circuits in Item (i) that is hard on average against  $\Sigma_{i-1}^p$ -oracle circuit. By Theorem 3.5, there is a function  $h_\ell : \{0, 1\}^\ell \rightarrow \{0, 1\}$  for some  $\ell = O(n^2)$  that is computable by  $\Pi_i\text{-SIZE}[\text{poly}(n) \cdot s_1(n)]$  circuits, such that for every  $\Sigma_{i-1}^p$ -oracle circuit  $D$  of size  $2^{\gamma \ell^{\gamma \delta}}$ ,

$$\Pr_{x \sim \{0,1\}^\ell} [h_\ell(x) = D(x)] \leq \frac{1}{2} + \frac{1}{2^{\gamma \ell^{\gamma \delta}}}.$$

By setting  $\varepsilon = (1/2) \cdot \delta \cdot \gamma$ , this contradicts Equation (14).

Therefore we conclude that  $\mathsf{T}_{\text{PV}}^i \not\vdash \mathsf{LB}^i(s_1, s_2, m, n_0)$ . □

## References

- [AB87] Noga Alon and Ravi B. Boppana. The monotone circuit complexity of boolean functions. *Combinatorica*, 7(1):1–22, 1987.
- [AB09] Sanjeev Arora and Boaz Barak. *Complexity Theory: A Modern Approach*. Cambridge University Press, 2009.
- [AK10] Eric Allender and Michal Koucký. Amplifying lower bounds by means of self-reducibility. *J. ACM*, 57(3):14:1–14:36, 2010.
- [And85] Alexander E. Andreev. On a method for obtaining lower bounds for the complexity of individual monotone functions. *Soviet Math. Dokl.*, 31(3):530–534, 1985.
- [AW09] Scott Aaronson and Avi Wigderson. Algebrization: A new barrier in complexity theory. *Transactions on Computation Theory (TOCT)*, 1(1), 2009.
- [Bey09] Olaf Beyersdorff. On the correspondence between arithmetic theories and propositional proof systems – a survey. *Mathematical Logic Quarterly*, 55(2):116–137, 2009.
- [BGS75] Theodore P. Baker, John Gill, and Robert Solovay. Relativizations of the  $P = ? NP$  Question. *SIAM J. Comput.*, 4(4):431–442, 1975.
- [BKKK20] Sam Buss, Valentine Kabanets, Antonina Kolokolova, and Michal Koucký. Expander construction in  $\text{VNC}^1$ . *Ann. Pure Appl. Log.*, 171(7):102796, 2020.
- [BKO20] Jan Bydzovsky, Jan Krajíček, and Igor C. Oliveira. Consistency of circuit lower bounds with bounded theories. *Logical Methods in Computer Science*, 16(2), 2020.
- [BM20] Jan Bydzovsky and Moritz Müller. Polynomial time ultrapowers and the consistency of circuit lower bounds. *Arch. Math. Log.*, 59(1-2):127–147, 2020.
- [Bus86] Samuel R. Buss. *Bounded Arithmetic*. Bibliopolis, 1986.
- [Bus95] Samuel R. Buss. Relating the bounded arithmetic and polynomial time hierarchies. *Ann. Pure Appl. Log.*, 75(1-2):67–77, 1995.
- [Bus97] Samuel R. Buss. Bounded arithmetic and propositional proof complexity. In *Logic of Computation*, pages 67–121. Springer Berlin Heidelberg, 1997.

- [Bus08] Samuel R. Buss. Bounded arithmetic, cryptography and complexity. *Theoria*, 63:147–167, 2008.
- [CHO<sup>+</sup>22] Lijie Chen, Shuichi Hirahara, Igor C. Oliveira, Ján Pich, Ninad Rajgopal, and Rahul Santhanam. Beyond natural proofs: Hardness magnification and locality. *J. ACM*, 69(4):25:1–25:49, 2022.
- [CJW19] Lijie Chen, Ce Jin, and Ryan Williams. Hardness magnification for all sparse NP languages. In *Symposium on Foundations of Computer Science (FOCS)*, pages 1240–1255, 2019.
- [CK07] Stephen A. Cook and Jan Krajíček. Consequences of the provability of  $NP \subseteq P/\text{poly}$ . *J. Symb. Log.*, 72(4):1353–1371, 2007.
- [CKKO21] Marco Carmosino, Valentine Kabanets, Antonina Kolokolova, and Igor C. Oliveira. Learn-uniform circuit lower bounds and provability in bounded arithmetic. In *Symposium on Foundations of Computer Science (FOCS)*, 2021.
- [CN10] Stephen A. Cook and Phuong Nguyen. *Logical Foundations of Proof Complexity*. Cambridge University Press, 2010.
- [Cob65] Alan Cobham. The intrinsic computational difficulty of functions. *Proc. Logic, Methodology and Philosophy of Science*, pages 24–30, 1965.
- [Coo75] Stephen A. Cook. Feasibly constructive proofs and the propositional calculus (preliminary version). In *Symposium on Theory of Computing (STOC)*, pages 83–97, 1975.
- [FGHK16] Magnus Gausdal Find, Alexander Golovnev, Edward A. Hirsch, and Alexander S. Kulikov. A better-than- $3n$  lower bound for the circuit complexity of an explicit function. In *Symposium on Foundations of Computer Science (FOCS)*, pages 89–98, 2016.
- [FLY22] Zhiyuan Fan, Jiayu Li, and Tianqi Yang. The exact complexity of pseudorandom functions and the black-box natural proof barrier for bootstrapping results in computational complexity. In *Symposium on Theory of Computing (STOC)*, pages 962–975, 2022.
- [Hås86] Johan Håstad. Almost optimal lower bounds for small depth circuits. In *Symposium on Theory of Computing (STOC)*, pages 6–20, 1986.
- [HVV06] Alexander Healy, Salil P. Vadhan, and Emanuele Viola. Using nondeterminism to amplify hardness. *SIAM J. Comput.*, 35(4):903–931, 2006.
- [Imp95] Russell Impagliazzo. Hard-core distributions for somewhat hard problems. In *Symposium on Foundations of Computer Science (FOCS)*, pages 538–545. IEEE Computer Society, 1995.
- [Jeř04] Emil Jeřábek. Dual weak pigeonhole principle, boolean complexity, and derandomization. *Ann. Pure Appl. Log.*, 129(1-3):1–37, 2004.
- [Jeř05] Emil Jeřábek. *Weak pigeonhole principle and randomized computation*. PhD thesis, 2005.
- [Jeř07a] Emil Jeřábek. Approximate counting in bounded arithmetic. *J. Symb. Log.*, 72(3):959–993, 2007.
- [Jeř07b] Emil Jeřábek. On independence of variants of the weak pigeonhole principle. *J. Log. Comput.*, 17(3):587–604, 2007.

- [Jer09] Emil Jerábek. Approximate counting by hashing in bounded arithmetic. *J. Symb. Log.*, 74(3):829–860, 2009.
- [Jer22] Emil Jerábek. Iterated multiplication in  $VTC^0$ . *Archive for Mathematical Logic*, pages 1–63, 2022.
- [KH82] Clement F. Kent and Bernard R. Hodgson. An arithmetical characterization of NP. *Theor. Comput. Sci.*, 21:255–267, 1982.
- [KKMP21] Robert Kleinberg, Oliver Korten, Daniel Mitropolsky, and Christos H. Papadimitriou. Total functions in the polynomial hierarchy. In *Innovations in Theoretical Computer Science Conference (ITCS)*, pages 44:1–44:18, 2021.
- [KO17] Jan Krajíček and Igor C. Oliveira. Unprovability of circuit upper bounds in Cook’s theory PV. *Logical Methods in Computer Science*, 13(1), 2017.
- [Koh08] Ulrich Kohlenbach. *Applied Proof Theory - Proof Interpretations and their Use in Mathematics*. Springer Monographs in Mathematics. Springer, 2008.
- [Kor21] Oliver Korten. The hardest explicit construction. In *Symposium on Foundations of Computer Science (FOCS)*, pages 433–444, 2021.
- [KPT91] Jan Krajíček, Pavel Pudlák, and Gaisi Takeuti. Bounded arithmetic and the polynomial hierarchy. *Ann. Pure Appl. Log.*, 52(1-2):143–153, 1991.
- [Kra92] Jan Krajíček. No counter-example interpretation and interactive computation. In Yiannis N. Moschovakis, editor, *Logic from Computer Science*, pages 287–293, New York, NY, 1992. Springer New York.
- [Kra95] Jan Krajíček. *Bounded Arithmetic, Propositional Logic, and Complexity Theory*. Encyclopedia of Mathematics and its Applications. Cambridge University Press, 1995.
- [Kra11] Jan Krajíček. On the proof complexity of the Nisan-Wigderson generator based on a hard  $NP \cap coNP$  function. *J. Math. Log.*, 11(1), 2011.
- [Kra19] Jan Krajíček. *Proof Complexity*. Encyclopedia of Mathematics and its Applications. Cambridge University Press, 2019.
- [Kra21] Jan Krajíček. Small circuits and dual weak PHP in the universal theory of p-time algorithms. *ACM Trans. Comput. Log.*, 22(2):11:1–11:4, 2021.
- [LC11] Dai Tri Man Le and Stephen A. Cook. Formalizing randomized matching algorithms. *Log. Methods Comput. Sci.*, 8(3), 2011.
- [LY22] Jiatu Li and Tianqi Yang.  $3.1n - o(n)$  circuit lower bounds for explicit functions. In *Symposium on Theory of Computing (STOC)*, pages 1180–1193, 2022.
- [Lê14] Dai Tri Man Lê. *Bounded Arithmetic and Formalizing Probabilistic Proofs*. PhD thesis, 2014.
- [MP20] Moritz Müller and Ján Pich. Feasibly constructive proofs of succinct weak circuit lower bounds. *Ann. Pure Appl. Log.*, 171(2), 2020.

- [Nis92] Noam Nisan. Pseudorandom generators for space-bounded computation. *Comb.*, 12(4):449–461, 1992.
- [NW94] Noam Nisan and Avi Wigderson. Hardness vs randomness. *J. Comput. Syst. Sci.*, 49(2):149–167, 1994.
- [Oja04] Kerry Ojakian. *Combinatorics in Bounded Arithmetic*. PhD thesis, 2004.
- [OS18] Igor C. Oliveira and Rahul Santhanam. Hardness magnification for natural problems. In *Symposium on Foundations of Computer Science (FOCS)*, pages 65–76, 2018.
- [Pap94] Christos H. Papadimitriou. *Computational Complexity*. Addison Wesley, 1994.
- [Pic14] Ján Pich. *Complexity Theory in Feasible Mathematics*. PhD thesis, 2014.
- [Pic15a] Ján Pich. Circuit lower bounds in bounded arithmetics. *Ann. Pure Appl. Log.*, 166(1):29–45, 2015.
- [Pic15b] Ján Pich. Logical strength of complexity theory and a formalization of the PCP theorem in bounded arithmetic. *Log. Methods Comput. Sci.*, 11(2), 2015.
- [PS21] Ján Pich and Rahul Santhanam. Strong co-nondeterministic lower bounds for NP cannot be proved feasibly. In *Symposium on Theory of Computing (STOC)*, 2021.
- [Raz85] Alexander A. Razborov. Lower bounds on the monotone complexity of some Boolean functions. *Doklady Akademii Nauk SSSR*, 281:798–801, 1985. English translation in: *Soviet Mathematics Doklady* 31:354–357, 1985.
- [Raz87] Alexander A. Razborov. Lower bounds on the size of constant-depth networks over a complete basis with logical addition. *Mathematicheskije Zametki*, 41(4):598–607, 1987.
- [Raz95a] Alexander A. Razborov. Bounded arithmetic and lower bounds in boolean complexity. In P. Clote and J. Remmel, editors, *Feasible Mathematics II*, pages 344—386. Birkhäuser, 1995.
- [Raz95b] Alexander A Razborov. Unprovability of lower bounds on circuit size in certain fragments of bounded arithmetic. *Izvestiya: mathematics*, 59(1):205, 1995.
- [RR97] Alexander A. Razborov and Steven Rudich. Natural proofs. *J. Comput. Syst. Sci.*, 55(1):24–35, 1997.
- [RSW22] Hanlin Ren, Rahul Santhanam, and Zhikun Wang. On the range avoidance problem for circuits. In *Symposium on Foundations of Computer Science (FOCS)*, 2022.
- [Smo87] Roman Smolensky. Algebraic methods in the theory of lower bounds for Boolean circuit complexity. In *Symposium on Theory of Computing (STOC)*, pages 77–82, 1987.
- [Sto76] Larry J. Stockmeyer. The polynomial-time hierarchy. *Theor. Comput. Sci.*, 3(1):1–22, 1976.
- [TC21] Iddo Tzameret and Stephen A. Cook. Uniform, integral, and feasible proofs for the determinant identities. *J. ACM*, 68(2):12:1–12:80, 2021.
- [TS00] A. S. Troelstra and H. Schwichtenberg. *Basic Proof Theory*. Cambridge Tracts in Theoretical Computer Science. Cambridge University Press, 2 edition, 2000.

- [Vad12] Salil P. Vadhan. Pseudorandomness. *Foundations and Trends® in Theoretical Computer Science*, 7(1–3):1–336, 2012.
- [Wil14] Ryan Williams. Nonuniform ACC circuit lower bounds. *J. ACM*, 61(1):2:1–2:32, 2014.
- [Wra76] Celia Wrathall. Complete sets and the polynomial-time hierarchy. *Theor. Comput. Sci.*, 3(1):23–33, 1976.
- [Zam96] Domenico Zambella. Notes on polynomially bounded arithmetic. *J. Symb. Log.*, 61(3):942–966, 1996.

## A Provability in $T_{PV}^i$

In this section, we further elaborate on the strength of the theories  $T_{PV}^i$ . Similarly to the relation between the complexity classes P, NP, and the different levels of PH, it is currently open if the theories  $T_{PV}^i$  form a proper hierarchy, i.e., if  $T_{PV}^j$  can prove more sentences than  $T_{PV}^i$  when  $j > i$ . However, as explained in this section, this is the case under standard computational hardness assumptions. Conversely, separating the theories would lead to new complexity class separations.

In Sections A.1 and A.2, we relate the relative strength of these theories to the hierarchy of total functions and to the polynomial time hierarchy, respectively. The results presented in these sections are closely related to results from [KPT91], which explore the strength of Buss’s theories  $S_2^i$  and  $T_2^i$  and related questions.

In Section A.3, we exhibit a complexity lower bound statement of comparatively higher quantifier complexity that is provable in  $T_{PV}^2$  (under a minimal assumption). In more detail, we show that if  $NP \not\subseteq (i.o.)P$  is true then it is provable in  $T_{PV}^2$ .

### A.1 Strength of $T_{PV}^i$ and the hierarchy of total functions

In this section, we show that separating the theories  $T_{PV}^i$  would lead to new complexity class separations. For instance, we prove that  $T_{PV}^2 = T_{PV}^1$  if and only if the search problem of every TFNP relation can be solved in polynomial time. A related result holds for every  $i \geq 1$  (see Theorem A.4 below for the precise statement).

The relationship between these theories and the corresponding complexity collapses provides evidence that the theories  $T_{PV}^i$  form a strict hierarchy. However, it also shows that unconditionally establishing that this is the case will be quite difficult.

For convenience, in the statements below we identify  $TF\Sigma_0^p$  with FP. We refer to Section 3.1 for definitions and to [KKMP21] for more information about total functions in the polynomial hierarchy. Abusing notation, in the statements below we view  $TF\Sigma_i^p$  as a class of search problems, i.e., given  $x$  the goal is to find  $y$  such that  $R(x, y)$  holds, where  $R \in TF\Sigma_i^p$ .

We will need the following lemmas, which can be proved using standard techniques from complexity and logic.

**Lemma A.1.** *For every  $i \geq 1$ ,  $P^{TF\Sigma_{i-1}^p} \subseteq \Sigma_{i-1}^p \subseteq P^{TF\Sigma_i^p}$ .*

**Lemma A.2.** *For every  $i \geq 1$ ,  $\Sigma_i^p \subseteq \Sigma_{i-1}^p$  if and only if  $\Sigma_i^p \subseteq P^{TF\Sigma_{i-1}^p}$ .*

**Lemma A.3.** *Let  $i \geq 1$ ,  $t(x)$  be an  $\mathcal{L}_{PV}$  term, and  $\phi(x, y)$  be a  $\Pi_{i-1}^b(\mathcal{L}_{PV})$ -formula. If  $T_{PV}^i \vdash \forall x \exists y \leq t(x) \phi(x, y)$ , then there exists a  $FP^{\Sigma_{i-1}^p}$  algorithm  $A(x)$  such that for every  $x \in \mathbb{N}$ ,  $\phi^{\mathbb{N}}(x, A(x))$  holds.*

The next theorem relates the relative strength of theories  $T_{PV}^i$  to the computational complexity of the search problems associated with the relations in  $TF\Sigma_j^p$ .

**Theorem A.4.** *For every  $i \geq 1$ , the following propositions hold:*

(i) *If  $TF\Sigma_i^p \subseteq FP^{TF\Sigma_{i-1}^p}$ , then  $T_{PV}^i \equiv T_{PV}^{i+1}$ .*

(ii) *If  $T_{PV}^i \equiv T_{PV}^{i+1}$ , then  $TF\Sigma_i^p \subseteq FP^{\Sigma_{i-1}^p}$ .*

*In particular,  $TFNP = FP$  if and only if  $T_{PV}^1 \equiv T_{PV}^2$ .*

*Proof.* (1) Assume that  $TF\Sigma_i^p \subseteq FP^{TF\Sigma_{i-1}^p}$ . We need to show that for every  $\varphi \in T_{PV}^{i+1}$ ,  $T_{PV}^i \vdash \varphi$ . Since it is enough to prove this for the axioms of  $T_{PV}^{i+1}$ , we can assume without loss of generality that  $\varphi = \forall x \exists y \leq t(x) \phi(x, y)$  for some  $\Pi_{i-1}^b$ -formula  $\phi$ . Let  $R \subseteq \{0, 1\}^* \times \{0, 1\}^*$  be the search problem such that  $(x, y) \in R$  if and only if  $y \leq t(x)$  and  $\phi(x, y)$ . Using the assumption in Item (i), we get that this search problem can be solved in  $FP^{TF\Sigma_{i-1}^p}$ . In particular, there is a  $\Sigma_{i-1}^b$ -formula  $\beta(x, y)$  that is total over  $\mathbb{N}$  and only accepts a pair  $(x, y)$  if  $(x, y) \in R$ . Thus  $T_{PV}^i \vdash \forall x \exists y \leq t(x) \beta(x, y)$  and  $T_{PV}^i \vdash \forall x \forall y \leq t(x) (\beta(x, y) \rightarrow \phi(x, y))$  by counting the quantifier complexity of these two sentences. It then follows that  $T_{PV}^i \vdash \varphi$ .

(2) Assume that  $T_{PV}^{i+1} \equiv T_{PV}^i$ . Let  $R \in TF\Sigma_i^p$  be a total relation such that for every  $(x, y) \in R$ ,  $|y| \leq |x|^c$ . Let  $\beta(x, y)$  be a  $\Pi_{i-1}^b$ -formula that captures over the standard model that  $(x, y) \in R$ . Then  $T_{PV}^{i+1} \vdash \forall x \exists y \in \{0, 1\}^{|x|^c} \beta(x, y)$ , which further means by the assumption in Item (ii) that  $T_{PV}^i \vdash \forall x \exists y \in \{0, 1\}^{|x|^c} \beta(x, y)$ . By Lemma A.3, we get that the search problem corresponding to  $R$  can be solved in  $FP^{\Sigma_{i-1}^p}$ .  $\square$

## A.2 Strength of $T_{PV}^i$ and the polynomial hierarchy

In this section, we relate the collapse of theories  $T_{PV}^i$  to a collapse of the polynomial hierarchy. More precisely, we show that if  $T_{PV}^{i+2} = T_{PV}^i$  then  $\Sigma_{i+1}^p = \Pi_{i+1}^p$ . Consequently, under the widely believed assumption that PH does not collapse, the theories  $T_{PV}^i$  can prove more sentences as  $i$  increases.

We will need a technical lemma from [KPT91] employed there to relate a certain collapse in Buss's hierarchy of theories of bounded arithmetic to a corresponding collapse of the polynomial hierarchy. First, we review the following statement.

**Principle  $\Omega(i)$ .** There is a constant  $k \in \mathbb{N}$  such that the following holds. For every relation  $P(x, y) \in \Pi_i^p$ , there are  $FP^{\Sigma_i^p}$  functions  $f_1(a), f_2(a, b_1), \dots, f_k(a, b_1, \dots, b_{k-1})$  such that:

- Either  $\forall z P^*(a, f_1(a), z)$  is true, or for every  $b_1$  s.t.  $\neg P^*(a, f_1(a), b_1)$ , it holds that:
- Either  $\forall z P^*(a, f_2(a, b_1), z)$  is true, or for every  $b_2$  s.t.  $\neg P^*(a, f_2(a, b_1), b_2)$ , it holds that:
- Either  $\forall z P^*(a, f_3(a, b_1, b_2), z)$  is true, or ...
- ...
- $\forall z P^*(a, f_k(a, b_1, b_2, \dots, b_k), z)$  is true;

where  $P^*(x, y, z) \triangleq |y| \leq |x| \wedge (y = 0 \vee P(x, y)) \wedge (|y| < |z| \leq |x| \rightarrow \neg P(x, z))$ .

**Lemma A.5** ([KPT91], Lemma 2.2). *For every  $i \geq 0$ , if Principle  $\Omega(i)$  is true, then  $\Sigma_{i+1}^p \subseteq P_{/poly}^{\Sigma_i^p}$  and thus also  $\Sigma_{i+2}^p = \Pi_{i+2}^p$ .*

**Theorem A.6.** *For every  $i \geq 1$ , if  $T_{PV}^i \equiv T_{PV}^{i+2}$ , then  $\Sigma_i^p \subseteq P_{/poly}^{\Sigma_{i-1}^p}$  and thus also  $\Sigma_{i+1}^p = \Pi_{i+1}^p$ .*



*Proof.* By Lemma A.5, it suffices to show that  $T_{PV}^i \equiv T_{PV}^{i+2}$  implies Principle  $\Omega(i-1)$ , for each  $i \geq 1$ . Assume that  $T_{PV}^i \equiv T_{PV}^{i+2}$ . For every relation  $P(x, y) \in \Pi_{i-1}^P$ , consider the  $\Pi_{i-1}^b(\mathcal{L}_{PV})$ -formula  $\alpha(x, y)$  that defines it and let  $\alpha^*(x, y, z)$  be defined as

$$\alpha^*(x, y, z) \triangleq |y| \leq |x| \wedge (y = 0 \vee \alpha(x, y)) \wedge (|y| < |z| \leq |x| \rightarrow \neg \alpha(x, z)).$$

Let  $\varphi \triangleq \forall x \exists |y| \leq |x| \forall |z| \leq |x| \alpha^*(x, y, z)$ . Since  $\mathbb{N} \models \varphi$  and  $\varphi$  is a  $\forall \Sigma_{i+1}^b$  sentence, we obtain that  $T_{PV}^{i+2} \vdash \varphi$  and thus  $T_{PV}^i \vdash \varphi$  by the assumption that  $T_{PV}^i \equiv T_{PV}^{i+2}$ . By Theorem 3.18, it follows that  $U_{PV}^i \vdash \varphi$ . Moreover, by Lemma 3.17, we know that

$$U_{PV}^i \vdash \forall x \exists y \forall z |y| \leq |x| \wedge (y = 0 \vee f_\alpha(x, y)) \wedge (|y| < |z| \leq |x| \rightarrow \neg f_\alpha(x, z)),$$

where  $f_\alpha^\mathbb{N}(x, y)$  is exactly  $P(x, y)$ . Principle  $\Omega(i-1)$  then follows directly from the KPT Witnessing Theorem (Theorem 3.11) and Theorem 3.21.  $\square$

### A.3 On the provability of $NP \not\subseteq (i.o.)P$

Recall that the axioms of  $T_{PV}^2$  consist of all true  $\forall \Sigma_1^b$ -sentences in the language  $\mathcal{L}_{PV}$ . In this section, we give a simple example of a complexity lower bound encoded by a collection of  $\forall \Sigma_2^b(\mathcal{L}_{PV})$ -sentences provable in  $T_{PV}^2$ , assuming the lower bound holds. (Note that the formalization below uses  $n \in \text{Log}$ , while our unprovability results hold even for  $n \in \text{LogLog}$ .)

**Theorem A.7.** *Assume that  $NP \not\subseteq (i.o.)P$ . For every polynomial-time Turing machine  $A$ , there is a constant  $n_0 \in \mathbb{N}$  such that  $T_{PV}^2$  proves*

$$\text{Fail}(A) \triangleq \forall n \in \text{Log} \exists \varphi(x_1, \dots, x_m) \in \{0, 1\}^n \left( n > n_0 \rightarrow \text{Error}(A, \varphi) \right),$$

where  $\varphi$  is an 3-CNF formula, and

$$\text{Error}(A, \varphi) \triangleq (\exists x \in \{0, 1\}^m \varphi(x) = 1 \wedge A(\varphi) = 0) \vee (\forall x \in \{0, 1\}^m \varphi(x) = 0 \wedge A(\varphi) = 1).$$

*Proof.* Assume that  $NP \not\subseteq (i.o.)P$ . Then  $3SAT \not\subseteq (i.o.)P$ , which means that for every polynomial-time Turing machine  $A$ , there exists a constant  $n_0$  such that  $A$  does not solve 3SAT on instances of length  $n > n_0$ . Let  $A$  be an arbitrary polynomial-time Turing machine. We would like to show that  $T_{PV}^2 \vdash \text{Fail}(A)$ .

**Search-to-Decision Reduction.** We firstly use a standard search-to-decision reduction to construct an efficient algorithm  $S$  that searches for a satisfying assignment, assuming that  $A$  solves SAT. An explicit description of  $S$  appears below (see Algorithm 11: Search-SAT Algorithm  $S$ ).

Without loss of generality, we assume that for every  $\varphi = \varphi_1 \in \{0, 1\}^n$  and  $i \in [m]$ , the corresponding formulas  $\varphi_i(x_i/0)$  and  $\varphi_i(x_i/1)$  can also be encoded as  $n$ -bit strings. Let  $A'$  be the following polynomial-time algorithm: given an instance  $\varphi \in \{0, 1\}^n$  encoding a 3-CNF formula; run  $S(\varphi) = (b, z)$ ; accept if and only if  $b = 2$  and  $\varphi(z) = 1$ .

**Claim A.8.** *There is a constant  $n_1 \in \mathbb{N}$  such that  $T_{PV}^2 \vdash \forall n \in \text{Log} \exists \varphi(x_1, \dots, x_m) \in \{0, 1\}^n \exists x \in \{0, 1\}^m (n > n_1 \rightarrow \varphi(x) = 1 \wedge A'(\varphi) = 0)$ .*

*Proof.* By the definition of  $A'$  we can see that it has only one-sided error, i.e., for every  $\varphi$  such that  $A'(\varphi) = 1$ ,  $\varphi$  is satisfiable. Since  $3SAT \not\subseteq (i.o.)P$ , the sentence is a  $\forall \Sigma_1^b$ -sentence that is true in the standard model provided that  $n_1$  is large enough, which further means that it is provable in  $T_{PV}^2$ .  $\square$

<p><b>Input :</b> A string <math>\varphi \in \{0, 1\}^n</math> encoding a 3-CNF formula.</p> <pre> 1 Let <math>\varphi_1(x_1, x_2, \dots, x_m)</math> be <math>\varphi</math>; 2 Let <math>z \in \{0, 1\}^m</math> be a string to be determined; 3 If <math>A(\varphi_1) = 0</math>, <b>return</b> <math>(0, 0)</math>; 4 <b>for</b> <math>i = 1, 2, \dots, m</math> <b>do</b> 5   <b>if</b> <math>A(\varphi_i(x_i/0)) = 1</math> <b>then</b> 6       Let <math>z_i = 0</math> and <math>\varphi_{i+1} = \varphi_i(x_i/0)</math>; 7   <b>else if</b> <math>A(\varphi_i(x_i/1)) = 1</math> <b>then</b> 8       Let <math>z_i = 1</math> and <math>\varphi_{i+1} = \varphi_i(x_i/1)</math>; 9   <b>else</b> 10      <b>return</b> <math>(1, \varphi_i)</math>; 11  <b>end</b> 12 <b>end</b> 13 <b>return</b> <math>(2, z)</math>; </pre>
--

**Algorithm 11:** Search-SAT Algorithm  $S$

**Claim A.9.** We have that  $T_{PV}^2 \vdash \forall n \in \text{Log } \forall \varphi(x_1, \dots, x_m) \in \{0, 1\}^n (A(\varphi) = 1 \wedge A'(\varphi) = 0 \rightarrow \exists \varphi' \in \{0, 1\}^n \text{Error}(A, \varphi'))$ .

*Proof.* Indeed, it is possible to establish even in PV that if  $\neg \text{Error}(A, \varphi')$  holds for every  $\varphi' \in \{0, 1\}^n$  then the search-to-decision reduction works as desired and consequently  $\neg(A(\varphi) = 1 \wedge A'(\varphi) = 0)$ . We omit the details.  $\square$

**Provability of the Hardness of 3SAT.** Now we prove in  $T_{PV}^2$  that  $\text{Fail}(A)$  holds for  $n_0 \triangleq n_1$ , where  $n_1 \in \mathbb{N}$  is the constant in Claim A.8. Arguing in the theory, let  $n \in \text{Log}$  be larger than  $n_1$ . Towards a contradiction, assume that for every  $\varphi(x_1, \dots, x_m) \in \{0, 1\}^n$ ,  $\neg \text{Error}(A, \varphi)$ . Let  $\varphi(x_1, \dots, x_m) \in \{0, 1\}^n$  be a 3-CNF formula from Claim A.8 such that  $\exists x \in \{0, 1\}^m (n > n_1 \rightarrow \varphi(x) = 1 \wedge A'(\varphi) = 0)$ . Since  $\varphi$  is satisfiable and by assumption  $\neg \text{Error}(A, \varphi)$ , we get that  $A(\varphi) = 1$ . Consequently, we have both  $A(\varphi) = 1 \wedge A'(\varphi) = 0$ . In turn, Claim A.9 yields the existence of  $\varphi' \in \{0, 1\}^n$  such that  $\text{Error}(A, \varphi')$ . This is in contradiction to the initial assumption on the correctness of  $A'$  on all instances of length  $n$ .  $\square$

## B Model-Theoretic Proof of the KPT Witnessing Theorem for $\forall \exists \forall \exists$ Sentences

**Theorem** (Theorem 3.13, restated). Let  $\mathcal{T}$  be a universal theory with vocabulary  $\mathcal{L}$ . Let  $\varphi$  be a quantifier-free  $\mathcal{L}$ -formula, and suppose that  $\mathcal{T} \vdash \forall x \exists y \forall z \exists w \varphi(x, y, z, w)$ . Then there is a constant  $\ell \geq 1$  and a sequence  $t_1, \dots, t_k$  of  $\mathcal{L}$ -terms such that

$$\mathcal{T} \vdash \forall x, z_1, \dots, z_k (\psi(z, t_1(z), z_1) \vee \psi(x, t_2(x, z_1), z_2) \vee \dots \vee \psi(x, t_k(z_1, \dots, z_{k-1}), z_k)),$$

where  $\psi(x, y, z) \triangleq \exists w \varphi(x, y, z, w)$ .

*Proof.* We verify that an argument described in [CN10, Section VIII.6] for the case of three quantifiers extends to sentences of this form. Let  $b, c_1, c_2, \dots$  be a list of new constant symbols, and let  $u_1, u_2, \dots$  be

an enumeration of all terms built from the functions and constants in  $\mathcal{L}$  together with  $b, c_1, c_2, \dots$ , where the only new constant symbols appearing in  $u_k$  are among  $b, c_1, \dots, c_{k-1}$ .

For convenience, let  $\psi(x, y, z) \triangleq \exists w \varphi(x, y, z, w)$ , as in the statement of the theorem. We will argue that there exists a constant  $k \geq 1$  such that no model of  $\mathcal{T}$  satisfies the sentence

$$\neg\psi(b, u_1, c_1) \wedge \neg\psi(b, u_2, c_2) \wedge \dots \wedge \neg\psi(b, u_k, c_k) .$$

This implies that every model of  $\mathcal{T}$  satisfies the negation of this sentence, and by the completeness theorem,

$$\mathcal{T} \vdash \psi(b, u_1, c_1) \vee \psi(b, u_2, c_2) \vee \dots \vee \psi(b, u_k, c_k) .$$

Since  $b, c_1, c_2, \dots$  are new constants and each term  $u_k$  depends only on  $b, c_1, \dots, c_{k-1}$  (among the new constant symbols), the result follows.

To show the remaining claim, we argue by contradiction. Suppose that no finite  $k$  satisfies the claim. Then, by compactness, we get that

$$\mathcal{T} \cup \{\neg\psi(b, u_1, c_1), \neg\psi(b, u_2, c_2), \neg\psi(b, u_3, c_3), \dots\}$$

admits a model  $\mathcal{M}$ . Consequently, using the definition of  $\psi$ ,

$$\mathcal{M} \models \mathcal{T} \cup \{\forall w \neg\varphi(b, u_1, c_1, w), \forall w \neg\varphi(b, u_2, c_2, w), \dots\}$$

Let  $\mathcal{T}^+ \triangleq \mathcal{T} \cup \{\forall w \neg\varphi(b, u_1, c_1, w), \forall w \neg\varphi(b, u_2, c_2, w), \dots\}$ . Since  $\mathcal{T}$  is a universal theory and  $\varphi$  is an open formula, it follows that  $\mathcal{T}^+$  is also a universal theory. For this reason, the substructure  $\mathcal{M}'$  of  $\mathcal{M}$  consisting of the denotations of the terms  $u_1, u_2, \dots$  (under  $\mathcal{M}$ ) is also a model of  $\mathcal{T}^+$ . Now it is not hard to prove that

$$\mathcal{M}' \models \mathcal{T} + \exists x \forall y \exists z \forall w \neg\varphi(x, y, z, w) ,$$

which contradicts the hypothesis of the theorem and completes the proof. To see this, it is enough to show that  $\mathcal{M}' \models \forall y \exists z \forall w \neg\varphi(b^{\mathcal{M}'}, y, z, w)$ . Given an arbitrary element  $m$  of  $\mathcal{M}'$ , by construction of  $\mathcal{M}'$ , there is some term  $u_k$  such that  $m = u_k^{\mathcal{M}'}(b^{\mathcal{M}'}, c_1^{\mathcal{M}'}, \dots, c_{k-1}^{\mathcal{M}'})$ . Since  $\mathcal{M}'$  is a model of  $\mathcal{T}^+$ , which includes the sentence  $\forall w \neg\varphi(b, u_k, c_k, w)$ , we get that  $\mathcal{M}' \models \forall w \neg\varphi(b^{\mathcal{M}'}, m, c_k^{\mathcal{M}'}, w)$ . This finishes the proof that  $\mathcal{M}' \models \forall y \exists z \forall w \neg\varphi(b^{\mathcal{M}'}, y, z, w)$ .  $\square$

We note that the argument described above does not extend to a larger number of quantifiers.

## C Self-Contained Proof of Theorem 4.20 via Herbrandization

The *no-counterexample interpretation* (see, e.g., [Koh08, Section 2.3] and [Kra92]) is a standard tool in proof theory to extract computational content from provable sentences of high quantifier complexity. In this section, we use this perspective to provide a different proof of Theorem 4.20. We refer to Section 4.7 for the necessary definitions and notation.

Let  $\mathcal{T}$  be a universal theory over  $\mathcal{L}$ , and let

$$\varphi(x) \triangleq \exists y_1 \forall x_1 \exists y_2 \dots \forall x_{k-1} \exists y_k \forall x_k \phi(x, \vec{x}, \vec{y})$$

be an  $\mathcal{L}$ -formula, where  $\phi$  is quantifier-free. The *Herbrand normal form* of  $\varphi(x)$  is defined as

$$\varphi^H(x) \triangleq \exists y_1 \exists y_2 \dots \exists y_k \phi(x, x_1/f_1(x, y_1), x_2/f_2(x, y_1, y_2), \dots, x_k/f_k(x, y_1, y_2, \dots, y_k), \vec{y}),$$

where  $f_1, f_2, \dots, f_k$  are new function symbols not in  $\mathcal{L}$ . By a simple model-theoretical argument,  $\mathcal{T} \vdash \forall x \varphi(x)$  if and only if  $\mathcal{T} \vdash \forall x \varphi^H(x)$ . Under the assumption that  $\mathcal{T} \vdash \forall x \varphi(x)$ , we can apply Theorem 3.10 to extract  $\mathcal{L}(f_1, f_2, \dots, f_k)$ -terms that witness the existential quantifiers. In particular, if  $\mathcal{T}$  is  $\mathsf{T}_{\text{PV}}$  and  $\mathcal{L}$  is  $\mathcal{L}_{\text{PV}}$ , this witnessing result implies that for every  $x \in \mathbb{N}$  and all interpretations of  $f_1, f_2, \dots, f_k$  over  $\mathbb{N}$ , we can find suitable  $y_1, y_2, \dots, y_k \in \mathbb{N}$  in polynomial-time given oracle access to  $f_1^{\mathbb{N}}, f_2^{\mathbb{N}}, \dots, f_k^{\mathbb{N}}$ .

Let  $\mathcal{M}$  be a structure over the vocabulary  $\mathcal{L}$  such that  $\mathcal{M} \models \mathcal{T}$  (e.g.,  $\mathcal{T} = \mathsf{T}_{\text{PV}}$  and  $\mathcal{M} = \mathbb{N}$ ), and let  $n_0$  be an object in the domain of  $\mathcal{M}$ . It is instructive to consider the following game on the board  $(\mathcal{M}, n_0)$ . There are two players in the game: a truthifier (or student) that claims  $\mathcal{M} \models \varphi(n_0)$ , and a falsifier (or teacher) that claims  $\mathcal{M} \models \neg \varphi(n_0)$ . In the  $i$ -th step, first the truthifier chooses an element  $n_i$  for  $y_i$ , then the falsifier chooses an element  $m_i$  for  $x_i$ . The truthifier (resp. falsifier) wins if and only if  $\mathcal{M} \models \varphi(n_0, m_1, \dots, m_k, n_1, \dots, n_k)$  holds (resp. does not hold). It is easy to see that  $\mathcal{M} \models \varphi(n_0)$  if and only if the truthifier has a winning strategy for the game on board  $(\mathcal{M}, n_0)$ . The interpretation of the function symbols  $f_1, \dots, f_k$  corresponds naturally to a strategy for the falsifier. The no-counterexample interpretation essentially means that if  $\mathcal{T} \vdash \forall x \varphi(x)$ , for every board  $(\mathcal{M}, n_0)$  and every strategy  $f_1, \dots, f_k$  of the falsifier, the truthifier has a winning strategy that can be expressed by terms in  $\mathcal{L}(f_1, f_2, \dots, f_k)$ . Next, we transform such a strategy into  $\mathcal{L}$ -strategies with ancillary information for the truthifier in the evaluation game of  $\varphi(x)$ .

**Theorem** (Reminder of Theorem 4.20). *Let  $\mathcal{T}$  be a universal theory over the language  $\mathcal{L}$  that is closed under if-then-else. Let  $\varphi(x)$  be the formula*

$$\begin{aligned} \varphi(x) \triangleq & \exists y_1 \leq t_1(x) \forall x_1 \leq s_1(x, y_1) \exists y_2 \leq t_2(x, y_1, x_1) \dots \forall x_{k-1} \leq s_{k-1}(x, y_1, x_1, \dots, y_{k-1}) \\ & \exists y_k \leq t_k(x, y_1, x_1, \dots, y_{k-1}, x_{k-1}) \forall x_k \leq s_k(x, y_1, x_1, \dots, y_k) \phi(x, x_1, \dots, x_k, y_1, \dots, y_k), \end{aligned}$$

where  $\phi(x, \vec{x}, \vec{y})$  is a quantifier-free  $\mathcal{L}$ -formula. If  $\mathcal{T} \vdash \forall x \varphi(x)$ , then there is a constant  $\ell \in \mathbb{N}$  and  $\mathcal{L}$ -strategies  $\tau_1^{\text{t}}, \tau_2^{\text{t}}, \dots, \tau_{\ell}^{\text{t}}$  (with ancillary information) such that, for any board  $(\mathcal{M}, n_0)$  and evaluation game of  $\varphi(x)$  on  $(\mathcal{M}, n_0)$ , for every strategy  $\tau^{\text{f}}$  of the falsifier:

- either  $\hat{\tau}_1^{\text{t}} \triangleq \tau_1^{\text{t}}[\emptyset]$  beats  $\tau^{\text{f}}$ ,
- or  $\hat{\tau}_2^{\text{t}} \triangleq \tau_2^{\text{t}}[\langle \hat{\tau}_1^{\text{t}} : \tau^{\text{f}} \rangle]$  beats  $\tau^{\text{f}}$ ,
- or  $\hat{\tau}_3^{\text{t}} \triangleq \tau_3^{\text{t}}[\langle \hat{\tau}_1^{\text{t}} : \tau^{\text{f}} \rangle, \langle \hat{\tau}_2^{\text{t}} : \tau^{\text{f}} \rangle]$  beats  $\tau^{\text{f}}$ ,
- ...,
- or  $\hat{\tau}_{\ell}^{\text{t}} \triangleq \tau_{\ell}^{\text{t}}[\langle \hat{\tau}_1^{\text{t}} : \tau^{\text{f}} \rangle, \langle \hat{\tau}_2^{\text{t}} : \tau^{\text{f}} \rangle, \dots, \langle \hat{\tau}_{\ell-1}^{\text{t}} : \tau^{\text{f}} \rangle]$  beats  $\tau^{\text{f}}$ .

*Proof.* We introduce Herbrandization functions  $f_1, f_2, \dots, f_k$  such that in  $\mathcal{L}^* \triangleq \mathcal{L} \cup \{f_1, \dots, f_k\}$ ,

$$\mathcal{T} \vdash \forall x \exists \vec{y} \leq \vec{t} \phi(x, \vec{x}^*, \vec{y}),$$

where  $x_j^* = f_j(x, y_1, y_2, \dots, y_j)$  for all  $j \in [k]$ . By Herbrand's Theorem (Theorem 3.10), there is a constant  $r \in \mathbb{N}$  and  $\mathcal{L}^*$ -terms  $q_j^i(x)$  ( $i \in [r], j \in [k]$ ) such that

$$\mathcal{T} \vdash \forall x \left( \bigvee_{i=1}^r \phi_i(x) \right),$$

where  $\phi_i(x) \triangleq \phi(x, x_1/f_1(x, q_1^i(x)), \dots, x_k/f_k(x, q_k^i(x)), y_1/q_1^i(x), \dots, y_k/q_k^i(x))$ .

We will translate  $(q_1^i, q_2^i, \dots, q_k^i)$  into  $\ell_i$   $\mathcal{L}$ -strategies  $\tau_{i,1}^{\text{t}}, \tau_{i,2}^{\text{t}}, \dots, \tau_{i,\ell_i}^{\text{t}}$  for some  $\ell_i \in \mathbb{N}$ , such that for every board  $(\mathcal{M} = (\mathcal{D}, \mathcal{I}), n_0)$  and every interpretation  $F_1, F_2, \dots, F_k$  of  $f_1, f_2, \dots, f_k$  over  $\mathcal{D}$  derived from

$\hat{\tau}^f$ , if  $\mathcal{M}(F_1, F_2, \dots, F_k) \models \phi_i(x/n_0)$ , then  $\tau_{i,1}^t, \tau_{i,2}^t, \dots, \tau_{i,\ell_i}^t$  will satisfy the conclusion of the theorem against the strategy  $\hat{\tau}^f$ . If this is possible, then

$$\tau_{1,1}^t, \tau_{1,2}^t, \dots, \tau_{1,\ell_1}^t, \tau_{2,1}^t, \tau_{2,2}^t, \dots, \tau_{2,\ell_2}^t, \dots, \tau_{r,1}^t, \tau_{r,2}^t, \dots, \tau_{r,\ell_r}^t$$

is a sequence of  $\mathcal{L}$ -strategies as required. The argument is as follows. Fix any board  $(\mathcal{M} = (\mathcal{D}, \mathcal{I}), n_0)$  and any strategy  $\tau^f$  of the falsifier. Let  $F_1, \dots, F_k$  be the interpretation of  $f_1, \dots, f_k$  corresponding to this strategy, i.e., for every  $j \in [k]$ ,

$$F_j(n, m_1, m_2, \dots, m_k) \triangleq \begin{cases} \text{the move of } \tau^f & \text{if } n = n_0 \text{ and } n_1, F_1(n_0, m_1), n_2, F_2(n_0, m_1, m_2), \\ \text{in the } j\text{-th step} & \dots, n_{j-1}, F_{j-1}(n_0, m_1, \dots, m_{j-1}) \text{ is a prefix of a} \\ 0 & \text{valid transcript over } (\mathcal{M}, n_0); \\ & \text{otherwise.} \end{cases}$$

Then there is an index  $i \in [r]$  such that  $\mathcal{M}(F_1, F_2, \dots, F_j) \models \phi_i(x/n_0)$  holds.

Before presenting the translation, we explain the main difficulty and how to address it. The issue is that  $(q_1^i, q_2^i, \dots, q_k^i)$  are  $\mathcal{L}^*$ -terms, while the desired strategy in the evaluation game consists of  $\mathcal{L}$ -terms only. For simplicity, suppose  $q_j^i(x)$  invokes a single function from the list  $f_1, \dots, f_k$  of new function symbols, and assume it is  $f_1(x, y_1)$ . The idea is to replace the computation  $F_1(w_1, w_2)$  over inputs  $w_1, w_2$  by forcing the falsifier to compute its value in a previously played game. To achieve this, we use that  $\tau^f$  is fixed. In other words, if  $w_1 = n_0$ , the falsifier must play and reveal  $F_1(n_0, w_1)$  if the truthifier plays  $w_1$  in the first round. (On the other hand, if  $w_1 \neq n_0$  we have  $F_1(w_1, w_2) = 0$  by definition.) Consequently, by playing more games we guarantee that the necessary information appears in the transcript, which allows us to replace calls to functions  $f_j$  and express the winning strategy using  $\mathcal{L}$ -terms. To streamline the presentation, in the description below we omit the trivial case where the first input to a function  $f_j$  is different than  $x$ , the input to the  $\mathcal{L}^*$ -terms  $q_j^i$  (corresponding to the case  $w_1 \neq n_0$  we have just explained).

Let  $\tau^f$  be the strategy specified by  $f_1, f_2, \dots, f_k$  (i.e.  $f_i$  denotes the falsifier's move in the  $i$ -th round). We prove by structural induction on the terms that we can decompose each  $q_j^i$  ( $j \in [k]$ ), which consists of  $\mathcal{L}$ -functions and  $f_1, f_2, \dots, f_k$ , into finitely many  $\mathcal{L}$ -strategies  $\tau_1^{i,j}, \tau_2^{i,j}, \dots, \tau_{d_{i,j}}^{i,j}$  and an  $\mathcal{L}$ -term  $p^{i,j}$  such that  $\mathcal{M}(F_1, F_2, \dots, F_j) \models q_j^i(n_0) = p^{i,j}(\Gamma(n_0))$  for every board  $(\mathcal{M}, \mathcal{I})$  and strategy  $\tau^f$  of the falsifier, where  $F_1, F_2, \dots, F_j$  is the interpretation of  $\tau^f$  corresponding to the strategy and  $\Gamma(n_0)$  is a sequence of transcripts produced as follows.

- For each  $u \in [d_{i,j}]$ , let  $\Gamma_u(n_0) \triangleq \langle \tau_u^{i,j}[\Gamma_1(n_0), \Gamma_2(n_0), \dots, \Gamma_{u-1}(n_0)] : \tau^f \rangle$ .
- Let  $\Gamma(n_0) \triangleq (\Gamma_1(n_0), \Gamma_2(n_0), \dots, \Gamma_{d_{i,j}}(n_0))$ .

Concretely, we translate each term as follows. Let the term be  $g(v_1, v_2, \dots, v_d)$ . By induction hypothesis, for every  $r \in [d]$ , we can decompose the term  $v_r$  into a sequence of  $c_r \in \mathbb{N}$   $\mathcal{L}$ -strategies  $\tau_1^r, \tau_2^r, \dots, \tau_{c_r}^r$  and an  $\mathcal{L}$ -term  $p_r$ , such that for every board  $(\mathcal{M}, m)$ ,  $\mathcal{M} \models v_r(n_0) = p_r(\Gamma^r(n_0))$  (where  $\Gamma^r(n_0)$  is the transcript of games as described above).

- If  $g(\cdot)$  is a function symbol in the original language  $\mathcal{L}$ , it is easy to see that the  $\mathcal{L}$ -term

$$g(p_1(\Gamma(n_0)), p_2(\Gamma(n_0)), \dots, p_d(\Gamma(n_0)))$$

and the strategies

$$(\tau_1^1, \tau_2^1, \dots, \tau_{c_1}^1, \tau_1^2, \tau_2^2, \dots, \tau_{c_2}^2, \dots, \tau_1^d, \tau_2^d, \dots, \tau_{c_d}^d)$$

provide what we want, where  $\Gamma(n_0) \triangleq (\Gamma^1(n_0), \Gamma^2(n_0), \dots, \Gamma^{c_d}(n_0))$ .

- If  $g(\cdot) = f_j$  for some  $j \in [k]$ , we define a new  $\mathcal{L}$ -strategy  $\tau^{f_j}$  as follows: suppose that the ancillary information consists of the transcripts  $\Gamma$  of  $\tau_1^1, \tau_2^1, \dots, \tau_{c_1}^1, \tau_1^2, \tau_2^2, \dots, \tau_{c_2}^2, \dots, \tau_1^d, \tau_2^d, \dots, \tau_{c_d}^d$  vs  $\tau^f$ ; in the  $i$ -th round for  $i \leq j$ , the truthifier's move is

$$\hat{p}_i(n_0, m_1, n_1, \dots, n_{i-1}, \Gamma) \triangleq \begin{cases} p_i(\Gamma) & p_i(\Gamma) \leq t_i(n_0, m_1, n_1, \dots, n_{i-1}) \\ 0 & \text{otherwise} \end{cases}$$

while in the remaining  $n - i$  rounds the truthifier always chooses 0. Note that  $\hat{p}_i$  is expressible since  $\mathcal{T}$  is closed under if-then-else. It is clear that the following  $\mathcal{L}$ -term  $v_g$  that takes  $\Gamma$  and the transcript  $\langle \tau^{f_j}[\Gamma] : \tau^f \rangle$  as input parameters outputs  $f_j(v_1(n_0), \dots, v_d(n_0))$ :

- If  $\bigvee_i p_i(\Gamma) > t_i(n_0, m_1, n_1, \dots, n_{i-1})$  holds, then  $v_g$  outputs 0.
- Otherwise,  $v_g$  outputs the  $j$ -th move of the falsifier in the transcript  $\langle \tau^{f_j}[\Gamma] : \tau^f \rangle$ .

Therefore, we can obtain a term  $v_g$  that simply reads the transcripts  $\Gamma, \langle \tau^{f_j}[\Gamma] : \tau^f \rangle$  and outputs  $f_j(v_1(n_0), \dots, v_d(n_0))$ , together with the strategies

$$\tau_1^1, \tau_2^1, \dots, \tau_{c_1}^1, \tau_1^2, \tau_2^2, \dots, \tau_{c_2}^2, \dots, \tau_1^d, \tau_2^d, \dots, \tau_{c_d}^d, \tau^{f_j},$$

as promised in the induction hypothesis.

Now we go back to the translation of  $(q_1^i, q_2^i, \dots, q_k^i)$  into  $\mathcal{L}$ -strategies. Assume that each  $q_j^i$  for  $j \in [k]$  has been decomposed into strategies  $\tau_1^{i,j}, \dots, \tau_{d_{i,j}}^{i,j}$  and a term  $p^{i,j}(\Gamma)$  as discussed above. Define an  $\mathcal{L}$ -strategy  $\tau^{q^i}$  with ancillary information as follows: suppose that the ancillary information is the transcripts  $\Gamma(n_0)$  of  $\tau^f$  vs

$$\tau_1^{i,1}, \tau_2^{i,1}, \dots, \tau_{d_{i,1}}^{i,1}, \tau_1^{i,2}, \tau_2^{i,2}, \dots, \tau_{d_{i,2}}^{i,2}, \dots, \tau_1^{i,k}, \tau_2^{i,k}, \dots, \tau_{d_{i,k}}^{i,k}$$

in which the latter strategies are given the transcripts of  $\tau^f$  vs previous strategies. In the  $j$ -th round, the truthifier's move is  $p^{i,j}(\Gamma(n_0))$ . By construction, it is easy to see that for every board  $(\mathcal{M}, n_0)$  and every strategy  $\tau^f$  of the falsifier, given correct ancillary information  $\Gamma(n_0)$ ,  $\tau^{q^i}$  will choose  $q_j^i(n_0)$  in the  $j$ -th round. Therefore, the strategy will beat  $\tau^f$  as long as  $\mathcal{M}(F_1, F_2, \dots, F_j) \models \phi_i(n_0)$ , where  $F_1, \dots, F_j$  constitute the interpretation of  $f_1, \dots, f_j$  corresponding to  $\tau^f$ . This completes the proof by previous discussions.  $\square$

## D Lemmas for Hardness Amplification

In this section, we discuss the technical lemmas stated without proof in Section 3.3. We follow the arguments in [HVV06]. For simplicity, we will only discuss below the parts that are different from [HVV06], referring to their paper for the omitted details.

**Reminder of Lemma 3.6.** *For every  $k \leq 2^n$ , there is an explicit computable generator  $G_k : \{0, 1\}^\ell \rightarrow (\{0, 1\}^n)^k$  that satisfies the requirements below:*

- (i) *There is an algorithm that computes the  $i$ -th block of  $G_k(\sigma)$  in  $\text{poly}(\ell, \log k)$  time given  $\sigma, i$ .*
- (ii)  *$G_k$  is indistinguishability-preserving for size  $t = k^2$ .*
- (iii)  *$G_k$  is  $2^{-n}$ -pseudorandom against branching programs of size  $2^n$  and block-size  $n$ .*

*Proof.* The only difference between this lemma and [HVV06, Lemma 5.12] is that in our definition, the indistinguishability-preserving property holds against  $\Sigma_{i-1}^p$ -oracle circuits instead of standard circuits, which will not cause any issue since their argument only requires mild closure properties of the adversary. For completeness, we sketch their proof here.

The generator  $G_k$  is defined as the XOR of two generators: a Nisan-Wigderson based generator  $NW_k : \{0, 1\}^{\ell_{NW}} \rightarrow (\{0, 1\}^n)^k$  that is efficiently computable and indistinguishability-preserving; and Nisan's unconditional PRG  $N_k : \{0, 1\}^{\ell_N} \rightarrow (\{0, 1\}^n)^k$  against (probabilistic) branching programs (see, e.g., [HVV06, Theorem 5.6] and [Nis92]). That is,  $G_k(x, y) \triangleq NW_k(x) \oplus N_k(y)$ . Both  $N_k$  and  $NW_k$  have seed length at most  $O(n^2)$ , hence  $G_k$  has seed length  $\ell = O(n^2)$ . Next, we discuss the properties of the generator.

- Both  $NW_k$  and  $N_k$  are efficiently computable in the sense that given  $\sigma$  and  $i$ , we can compute the  $i$ -th block of the output in  $\text{poly}(\ell, \log k)$  time. Therefore Item (i) holds.
- To prove Item (ii), we need to show that any indistinguishability-preserving generator XORed with a fixed string is still indistinguishability-preserving. Towards a contradiction, assume that  $G_k$  is not indistinguishability-preserving. This means that there are  $f_1, \dots, f_k, g_1, \dots, g_k$  such that for every  $i \in [k]$ ,  $x \| f_i(x) \approx_\varepsilon^s x \| g_i(x)$  for  $x \sim \{0, 1\}^n$ , while for  $(\sigma_1, \sigma_2) \sim \{0, 1\}^{\ell_{NW}} \times \{0, 1\}^{\ell_N}$  and  $(X_1, \dots, X_k) \triangleq NW_k(\sigma_1) \oplus N_k(\sigma_2)$ ,

$$\sigma_1 \| \sigma_2 \| f_1(X_1) \| \dots \| f_k(X_k) \not\approx_{\varepsilon}^{s-t} \sigma_1 \| \sigma_2 \| g_1(X_1) \| \dots \| g_k(X_k).$$

By an averaging argument, there is a  $\sigma_2^* \in \{0, 1\}^{\ell_N}$  such that

$$\sigma_1 \| f_1(X_1 \oplus y_1) \| \dots \| f_k(X_k \oplus y_k) \not\approx_{\varepsilon}^{s-t} \sigma_1 \| g_1(X_1 \oplus y_1) \| \dots \| g_k(X_k \oplus y_k), \quad (15)$$

where  $(y_1, \dots, y_k) \triangleq N_k(\sigma_2^*)$ . Let  $f'_i(x) \triangleq f_i(x \oplus y_i)$  and  $g'_i(x) \triangleq g_i(x \oplus y_i)$  for  $i \in [k]$ . Clearly for every  $i \in [k]$ ,  $x \| f'_i(x) \approx_\varepsilon^s x \| g'_i(x)$ ,<sup>36</sup> which is impossible since  $NW_k$  is indistinguishability-preserving but Equation (15) holds.

- Similarly, we can show that since  $N_k$  is  $2^{-n}$ -pseudorandom against branching programs of size  $2^n$ , after XORed with another generator,  $G_k$  is still  $2^{-n}$ -pseudorandom against branching programs of size  $2^n$ . This implies Item (iii).

It remains to verify that the Nisan-Wigderson based generator  $NW_k$  is indistinguishability-preserving for size  $k^2$  against  $\Sigma_{i-1}^p$ -oracle circuits. Let  $\ell = O(n^2)$  and  $S_1, S_2, \dots, S_k \subseteq [\ell]$  be an  $(\ell, n, \log k)$ -design (see Section 3.2 and [Nis92]). Then  $NW_k : \{0, 1\}^\ell \rightarrow (\{0, 1\}^n)^k$  is defined as

$$NW_k(\sigma) \triangleq (\sigma|_{S_1}, \sigma|_{S_2}, \dots, \sigma|_{S_k}).$$

Let  $f_1, \dots, f_k, g_1, \dots, g_k$  be probabilistic functions such that for every  $i \in [k]$ ,  $x \| f_i(x) \approx_\varepsilon^s x \| g_i(x)$  for  $x \sim \{0, 1\}^n$ . Suppose, for the sake of contradiction, that

$$\sigma \| f_1(\sigma|_{S_1}) \| \dots \| f_k(\sigma|_{S_k}) \not\approx_{k \cdot \varepsilon}^{s-k^2} \sigma \| g_1(\sigma|_{S_1}) \| \dots \| g_k(\sigma|_{S_k}). \quad (16)$$

For every  $i \in [0, k]$ , we define the hybrid distribution

$$H_i = \sigma \| g_1(\sigma|_{S_1}) \| \dots \| g_i(\sigma|_{S_i}) \| f_{i+1}(\sigma|_{S_{i+1}}) \| \dots \| f_k(\sigma|_{S_k}).$$

<sup>36</sup>There is no loss in the circuit size of the adversary if we define the circuit model so that NOT gates are free.



Then the distinguisher  $D$  for Equation (16), which is a  $\Sigma_{i-1}^p$ -oracle circuit of size  $s - k^2$ , can distinguish between  $H_i$  and  $H_{i+1}$  with advantage at least  $\varepsilon$  for some  $i \in [0, k-1]$ . Note that

$$\begin{aligned} H_i &= \sigma \|g_1(\sigma|_{S_1})\| \dots \|g_i(\sigma|_{S_i})\| f_{i+1}(\sigma|_{S_{i+1}}) \|f_{i+2}(\sigma|_{S_{i+2}})\| \dots \|f_k(\sigma|_{S_k})\| \\ H_{i+1} &= \sigma \|g_1(\sigma|_{S_1})\| \dots \|g_i(\sigma|_{S_i})\| g_{i+1}(\sigma|_{S_{i+1}}) \|f_{i+2}(\sigma|_{S_{i+2}})\| \dots \|f_k(\sigma|_{S_k})\| \end{aligned} \quad \text{and}$$

differ only on the  $(i+2)$ -th part:  $H_i$  has  $f_{i+1}(\sigma|_{S_{i+1}})$  while  $H_{i+1}$  has  $g_{i+1}(\sigma|_{S_{i+1}})$ .

By an averaging argument, we can fix all the bits of  $\sigma$  outside of  $S_{i+1}$  so that  $\hat{H}_i$  and  $\hat{H}_{i+1}$  are still distinguishable with advantage  $\varepsilon$ , where  $\hat{H}_i$  and  $\hat{H}_{i+1}$  refer to the distribution  $H_i$  and  $H_{i+1}$  after we fix the bits of  $\sigma$  outside of  $S_{i+1}$ . Since for every  $j \neq i+1$ ,  $|S_j \cap S_{i+1}| \leq \log k$ , we can construct a  $\Sigma_{i-1}^p$ -oracle circuit of size at most  $(s - k^2) + 2^{\log k} \cdot k = s$  that hardwires all possibilities for the common parts of  $H_i$  and  $H_{i+1}$  such that:

- Given the unfixed bits of  $\sigma$  and  $f_{i+1}(\sigma)$ , it generates  $\hat{H}_i$  and outputs  $D(\hat{H}_i)$ .
- Given the unfixed bits of  $\sigma$  and  $g_{i+1}(\sigma)$ , it generates  $\hat{H}_{i+1}$  and outputs  $D(\hat{H}_{i+1})$ .

Since  $D$  can distinguish between  $\hat{H}_i$  and  $\hat{H}_{i+1}$  with advantage  $\varepsilon$ , the circuit above can distinguish between  $\sigma \|f_{i+1}(\sigma|_{S_{i+1}})$  and  $\sigma \|g_{i+1}(\sigma|_{S_{i+1}})$  with advantage  $\varepsilon$ . This leads to a contradiction.  $\square$

**Reminder of Lemma 3.7.** For every  $i \geq 1$ ,  $\delta(n) = 1/\text{poly}(n)$ , and  $k = k(n)$  such that  $n^{\omega(1)} \leq k \leq 2^n$ , there is a function  $C_k : \{0, 1\}^k \rightarrow \{0, 1\}$  such that:

- (i)  $\text{NoiseStab}_\delta[C_k] \leq 1/k^{\Omega(1)}$ ;
- (ii) For every  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  computable by  $\Sigma_i\text{-SIZE}[s(n)]$  circuits,  $(C_k \circ f^{\otimes k}) \circ G_k : \{0, 1\}^\ell \rightarrow \{0, 1\}$  is computable by  $\Sigma_i\text{-SIZE}[\text{poly}(n) \cdot s(n)]$  circuits.
- (iii)  $C_k$  is computable by a branching program of size  $\text{poly}(n) \cdot k$  and also by a deterministic circuit of size  $\text{poly}(n) \cdot k$ .

*Proof.* Let  $\delta = \delta(n) \geq 1/\text{poly}(n)$  and  $k = k(n)$  such that  $n^{\omega(1)} \leq k \leq 2^n$ . We will define  $C_k$  as the composition of two functions defined as follows:

- The *recursive-majority function*  $\text{RMaj}_r : \{0, 1\}^{3^r} \rightarrow \{0, 1\}$  is recursively defined by

$$\begin{aligned} \text{RMaj}_1(x_1, x_2, x_3) &\triangleq \text{Maj}(x_1, x_2, x_3) \\ \text{RMaj}_r(x_1, \dots, x_{3^r}) &\triangleq \text{RMaj}_{r-1}(\text{Maj}(x_1, x_2, x_3), \dots, \text{Maj}(x_{3^{r-2}}, x_{3^{r-1}}, x_{3^r})) \end{aligned}$$

where  $\text{Maj}(x_1, x_2, x_3)$  is the majority value among  $x_1, x_2, x_3 \in \{0, 1\}$ .

- The *tribes function* of  $k$  bits is defined by

$$\text{Tribes}_k(x_1, \dots, x_k) \triangleq (x_1 \wedge \dots \wedge x_b) \vee (x_{b+1} \wedge \dots \wedge x_{2b}) \vee \dots \vee (x_{k-b+1} \wedge \dots \wedge x_k),$$

where  $b = O(\log k)$  is the largest integer such that  $(1 - 2^{-b})^{k/b} \geq 1/2$ .

Let  $r \triangleq c \cdot \log(1/\delta)$  for a constant  $c$  to be determined later. Assuming without loss of generality that  $r$  and  $k/3^r$  are integers, we define  $C_k : \{0, 1\}^k \rightarrow \{0, 1\}$  by

$$C_k \triangleq \text{Tribes}_{k/3^r} \circ \text{RMaj}_r^{\otimes k/3^r}.$$

As [HVV06, Section 5.5] in the proof of Lemma 5.15, we know that for some sufficiently large constant  $c$ , the noise stability of  $C_k$  is at most  $1/k^{\Omega(1)}$ . Also they showed that  $C_k$  can be computed by a branching program of size  $\text{poly}(n) \cdot k$  and a deterministic circuit of size  $\text{poly}(n) \cdot k$ .

It remains to determine the complexity of  $(C_k \circ f^{\otimes k}) \circ G_k$  for  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  computable by  $\Sigma_i\text{-SIZE}[s(n)]$  circuits. Consider the following  $\Sigma_i$ -circuit. We first guess (using non-determinism) a clause  $K$  of the upper  $\text{Tribes}_{k/3^r}$  function that is satisfied. For every  $\text{RMaj}_r$  function feeding into this clause (there are  $b = O(\log k) = \text{poly}(n)$  such  $\text{RMaj}_r$  functions), we guess the input bits of the upper  $C_k$  sub-circuit (or equivalently, the output bits of the lower  $f$  functions) that are 1 and

- (i) we verify that these input bits that are 1 make the clause  $K$  accept, which can be done by a deterministic circuit of size  $\text{poly}(3^r) = \text{poly}(n)$  since  $\text{RMaj}$  is a monotone function;
- (ii) for every guessed input bit of  $C_k$  (or equivalently, the output bit of one of  $f$  in the middle  $f^{\otimes k}$  layer) that is supposed to be 1, we use the  $\Sigma_i\text{-SIZE}[s(n)]$  circuit for  $f$  to verify that it is indeed 1. The input to this function  $f$  is one of the  $n$ -bit blocks of the output of  $G_k$ , which can be computed by a deterministic algorithm in  $\text{poly}(\ell, \log k) = \text{poly}(n)$  time (see Lemma 3.6).

The overall  $\Sigma_i$ -circuit complexity of  $(C_k \circ f^{\otimes k}) \circ G_k$  is at most  $\text{poly}(n) \cdot s(n)$ .  $\square$

**Reminder of Lemma 3.9.** Assume that  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  is  $\delta$ -hard for size  $s = n^{\omega(1)}$ . There is a  $\delta'$ -random function  $g$  with  $\delta' \in [\delta/2, \delta]$  such that the  $\text{Amp}_f(\sigma) : \{0, 1\}^\ell \rightarrow \{0, 1\}$  has hardness

$$\frac{1}{2} - \frac{\text{ExpBias}[(C \circ g^{\otimes k}) \circ G]}{2} - \frac{k}{s^{1/3}}$$

for size  $\Omega(s^{1/3} / \log(s/\delta)) - k^2 - \text{poly}(n) \cdot k$ .

Before proving this lemma, we need to verify that Impagliazzo's hardcore lemma (see, e.g., [AB09, Section 19.1.2]) holds against adversaries with access to  $\Sigma_{i-1}^p$  oracles.

**Lemma D.1** (Generalized version of Impagliazzo's Hardcore Lemma). Assume that  $2n < s < 0.001 \cdot (\varepsilon\delta)^2 \cdot 2^n/n$ . Let  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  be a balanced function that is  $\delta$ -hard for  $\Sigma_{i-1}^p$ -oracle circuits of size  $s$ . There exists a  $\delta'$ -random function  $g : \{0, 1\}^n \rightarrow \{0, 1\}$  such that  $X \| f(X) \approx_{\varepsilon}^{s'} X \| g(X)$  for  $X \sim \{0, 1\}^n$ , where  $s' = \Omega(s\varepsilon^2 / \log(1/(\delta\varepsilon)))$  and  $\delta' \in [\delta/2, \delta]$ .

*Proof Sketch.* We follow the proof presented in [AB09, Section 19.1.2] based on the *min-max theorem* for zero-sum games (also see, e.g., [Imp95]). We say that a distribution  $\mathcal{H}$  over  $\{0, 1\}^n$  has density  $\delta$  if for every  $x \in \{0, 1\}^n$ ,  $\mathcal{H}(x) \leq 1/(\delta 2^n)$ . Let  $\delta_1 = 0.99\delta$ . We first show that there is a distribution  $\mathcal{H}$  of density  $\delta_1$  such that for every  $\Sigma_{i-1}^p$ -oracle circuit  $C$  of size  $s'$ ,  $\Pr[f(x) = C(x)] < 1/2 + \varepsilon/2$  for  $x \sim \mathcal{H}$ .

Towards a contradiction, we assume that such distribution does not exist. By a game-theoretic argument using the min-max theorem, we can construct a distribution  $\mathcal{C}$  over  $\Sigma_{i-1}^p$ -oracle circuits of size  $s'$  such that for every distribution  $\mathcal{H}$  of density  $\delta_1$ , a random  $C \sim \mathcal{C}$  can approximate  $f$  over  $\mathcal{H}$  with error  $\leq 1/2 - \varepsilon/2$ .

An input  $x \in \{0, 1\}^n$  is said to be *bad* if  $\Pr[C(x) \neq f(x)] > 1/2 - \varepsilon/2$  for  $C \sim \mathcal{C}$ . It is said to be *good* otherwise. There are at most  $\delta_1 \cdot 2^n$  bad inputs, since otherwise we can let  $\mathcal{H}$  be the uniform distribution over a set of  $\delta_1 \cdot 2^n$  bad inputs and violate the aforementioned property of  $\mathcal{C}$ . Let  $t = O(\varepsilon^{-2} \log(1/(\delta\varepsilon)))$  and  $C$  be the following probabilistic circuit (with  $\Sigma_{i-1}^p$  oracles): given input  $x$ , obtain  $t$  independent samples  $C_1, \dots, C_t \sim \mathcal{C}$ , and output the majority of  $C_1(x), \dots, C_t(x)$ . This probabilistic circuit has size at most  $t \cdot s' \leq s$ . By the Chernoff bound, it computes  $f(x)$  for any good  $x$  with error at most  $\exp(-\Omega(\varepsilon^2 t)) \leq 0.001 \cdot \delta$ . This means that for a uniformly random  $x \sim \{0, 1\}^n$ , the probabilistic  $\Sigma_{i-1}^p$ -oracle circuit (and also deterministic  $\Sigma_{i-1}^p$ -oracle circuit by an averaging argument) can approximate  $f(x)$  with error at most  $\delta_1 + \delta/2 \leq \delta$  for an  $x \sim \{0, 1\}^n$ , which is impossible.

We then prove via a probabilistic argument that there is a subset  $H$  of size  $\delta' \in [\delta/2, \delta]$  such that no  $\Sigma_{i-1}^p$ -oracle circuit of size  $s$  can approximate  $f$  on  $H$  with advantage  $\varepsilon$ . Let  $H$  be a *random* subset defined

as follows: for every  $x \in \{0, 1\}^n$ , we let  $x \in H$  independently with probability  $\mathcal{H}(x)$ . By a “concentration bound then union bound” argument, we get with non-zero probability that  $H$  has size  $\delta' \in [\delta/2, \delta]$  and for every  $C$  of size  $s$ ,  $\Pr[f(x) = C(x)] \leq 1/2 + 1/\varepsilon$ . This means that the  $\delta'$ -random function  $g$  defined over  $H$  satisfies the conditions of the lemma.  $\square$

*Proof of Lemma 3.9.* Assume that  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  is  $\delta$ -hard for size  $s = n^{\omega(1)}$ . By Impagliazzo’s hardcore lemma, there is a  $\delta'$ -random function  $g : \{0, 1\}^n \rightarrow \{0, 1\}$  such that  $X \| f(X) \approx_{\varepsilon}^{s'} X \| g(X)$  for  $X \sim \{0, 1\}^n$ , where  $s' = \Omega(s\varepsilon^2 / \log(1/(\delta\varepsilon)))$  and  $\delta' \in [\delta/2, \delta]$ . Since  $G$  is indistinguishability-preserving for size  $k^2$ , we get that

$$\sigma \| f(X_1) \| \dots \| f(X_k) \approx_{k\varepsilon}^{s' - k^2} \sigma \| g(X_1) \| \dots \| g(X_k),$$

where  $\sigma \sim \{0, 1\}^\ell$  and  $(X_1, \dots, X_k) = G(\sigma)$ . Since  $C_k$  has complexity bounded by  $\text{poly}(n) \cdot k$  this further means that

$$\sigma \| C_k(f(X_1), \dots, f(X_k)) \approx_{k\varepsilon}^{s''} \sigma \| C_k(g(X_1), \dots, g(X_k)),$$

where  $s'' = s' - k^2 - \text{poly}(n) \cdot k$ . Note that

$$C_k(f(X_1), \dots, f(X_k)) = (C_k \circ f^{\otimes k}) \circ G_k(\sigma) \quad \text{and} \quad C_k(g(X_1), \dots, g(X_k)) = (C_k \circ g^{\otimes k}) \circ G_k(\sigma).$$

Also we can see that for every probabilistic function  $h$ , the statistical distance between  $X \| h(X)$  and  $X \| b$  for  $X \sim \{0, 1\}^n$  and  $b \sim \{0, 1\}$  is exactly  $\text{ExpBias}[h]/2$  (see, e.g., [HVV06, Lemma 3.4]). Therefore we know that

$$\Delta(\sigma \| (C_k \circ g^{\otimes k}) \circ G_k(\sigma), \sigma \| b) \leq \frac{\text{ExpBias}[(C_k \circ g^{\otimes k}) \circ G_k]}{2},$$

where  $\sigma \sim \{0, 1\}^\ell$  and  $b \sim \{0, 1\}$ . This further means that  $\sigma \| (C_k \circ f^{\otimes k}) \circ G_k(\sigma)$  and  $\sigma \| b$  are  $k\varepsilon + (1/2) \cdot \text{ExpBias}[(C_k \circ g^{\otimes k}) \circ G_k]$  indistinguishable for size  $s''$ . By setting  $\varepsilon = s^{-1/3}$ , we obtain the lemma.  $\square$

## E The Counting Lemma: Existence of a Good Restriction

**Notation.** Recall that for  $m \geq 1$ , a set  $S \subseteq \{0, 1\}^{[m]}$ , and a string  $a \in \{0, 1\}^I$ , where  $I \subseteq [m]$ , we define the *restriction of  $S$  with respect to  $a$*  as the set

$$S \upharpoonright_a \triangleq \{w \in S \mid w|_I = a\}.$$

For a non-empty set  $U$  and a set  $S \subseteq U$ , we define  $\text{dens}_U(S) \triangleq |S|/|U|$ .

For simplicity of the exposition, we consider without loss of generality restrictions with respect to the first  $m_1$  input bits. Let  $S \subseteq \{0, 1\}^m$ , where  $m = m_1 + m_2$ . Suppose that  $\text{dens}_{\{0, 1\}^m}(S) = \delta$ . Now let  $T \subseteq S$  be a set such that  $\text{dens}_S(T) > 2/3$ . The following result appears implicit in [Kra11, Pic15a].

**Lemma E.1** (Counting Lemma). *Under these assumptions, there is  $a \in \{0, 1\}^{m_1}$  such that*

$$\frac{|S \upharpoonright_a|}{2^{m_2}} \geq \frac{1}{100} \cdot \delta \quad \text{and} \quad \frac{|T \upharpoonright_a|}{|S \upharpoonright_a|} \geq \frac{2}{3} - \frac{1}{100}. \quad (17)$$

*Proof.* Suppose this is not the case, i.e., for every  $a \in \{0, 1\}^{m_1}$ , at least one of the two inequalities above does not hold. We use this to contradict  $|T| > (2/3) \cdot |S| = (2/3) \cdot \delta \cdot 2^m$ . Under the assumption, and using that  $T \restriction_a \subseteq S \restriction_a$ ,

$$\begin{aligned}
|T| &= \sum_{a \in \{0,1\}^{m_1}} |T \restriction_a| \\
&\leq \sum_{a \in \{0,1\}^{m_1}} \left( \frac{1}{100} \cdot \delta \cdot 2^{m_2} + \left( \frac{2}{3} - \frac{1}{100} \right) |S \restriction_a| \right) \\
&= 2^{m_1+m_2} \cdot \frac{1}{100} \cdot \delta + \left( \frac{2}{3} - \frac{1}{100} \right) \cdot |S| \\
&= \frac{\delta}{100} \cdot 2^m + \left( \frac{2}{3} - \frac{1}{100} \right) \cdot \delta \cdot 2^m \\
&= \frac{2}{3} \cdot \delta \cdot 2^m.
\end{aligned} \tag{18}$$

This completes the proof. □