

# Provability of the Circuit Size Hierarchy and Its Consequences

Marco Carmosino\*

Valentine Kabanets†

Antonina Kolokolova‡

Igor C. Oliveira§

Dimitrios Tsintsilidas¶

July 30, 2024

## Abstract

The *Circuit Size Hierarchy* ( $\text{CSH}_b^a$ ) states that if  $a > b \geq 1$  then the set of functions on  $n$  variables computed by Boolean circuits of size  $n^a$  is strictly larger than the set of functions computed by circuits of size  $n^b$ . This result, which is a cornerstone of circuit complexity theory, follows from the *non-constructive* proof of the existence of functions of large circuit complexity obtained by Shannon in 1949.

Are there more “constructive” proofs of the Circuit Size Hierarchy? Can we quantify this? Motivated by these questions, we investigate the provability of  $\text{CSH}_b^a$  in theories of bounded arithmetic. Among other contributions, we establish the following results:

(i) Given any  $b > 1$ ,  $\text{CSH}_b^a$  is provable in Buss’s theory  $T_2^2$  for  $a > b + 1$ .

(ii) In contrast, if there are constants  $a > b > 1$  such that  $\text{CSH}_b^a$  is provable in the theory  $T_2^1$ , then there is a constant  $\varepsilon > 0$  such that  $\text{P}^{\text{NP}}$  requires non-uniform circuits of size  $n^{1+\varepsilon}$ .

In other words, an improved *upper bound* on the proof complexity of  $\text{CSH}_b^a$  would lead to new *lower bounds* in complexity theory.

We complement these results with a proof of the *Formula Size Hierarchy* ( $\text{FSH}_b^a$ ) in  $\text{PV}_1$  with parameters  $a > 2$  and  $b = 3/2$ . This is in contrast with typical formalizations of complexity lower bounds in bounded arithmetic, which require  $\text{APC}_1$  or stronger theories and are not known to hold even in  $T_2^1$ .

---

\*IBM Research, USA. E-mail: marco@ntime.org

†Simon Fraser University, Canada. E-mail: kabanets@cs.sfu.ca

‡Memorial University of Newfoundland, Canada. E-mail: kol@mun.ca

§University of Warwick, UK. E-mail: igor.oliveira@warwick.ac.uk

¶University of Warwick, UK. E-mail: dimitrios.tsintsilidas@warwick.ac.uk

21 **Contents**

22	<b>1 Introduction</b>	<b>3</b>
23	1.1 Context and Motivation . . . . .	3
24	1.2 Results . . . . .	3
25	1.3 Techniques . . . . .	5
26	<b>2 Preliminaries</b>	<b>6</b>
27	2.1 Complexity Theory . . . . .	6
28	2.2 Bounded Arithmetic . . . . .	7
29	2.2.1 Logical Theories . . . . .	7
30	2.2.2 The KPT Witnessing Theorem . . . . .	8
31	<b>3 Circuit Size Hierarchies in Bounded Arithmetic</b>	<b>9</b>
32	3.1 Explicit Circuit Lower Bounds from Provability in $PV_1$ and $T_2^1$ . . . . .	9
33	3.2 Extracting All the Hardness from Proofs of a Succinct Hierarchy Theorem . . . . .	11
34	3.3 Formalization in $T_2^2$ . . . . .	12
35	3.4 On the Gap Between $T_2^1$ and $T_2^2$ . . . . .	13
36	<b>4 Provability of Formula Size Bounds in <math>PV_1</math></b>	<b>14</b>
37	4.1 Subbotovskaya’s Lower Bound . . . . .	14
38	4.1.1 High-Level Details of the Formalization . . . . .	14
39	4.1.2 On the Low-Level Details of the Formalization . . . . .	19
40	4.2 Upper Bound . . . . .	20
41	4.3 Formula Size Hierarchy . . . . .	23
42	<b>A Proof of the KPT Theorem for <math>\forall E \exists E</math> Sentences</b>	<b>25</b>

# 1 Introduction

## 1.1 Context and Motivation

The existence of Boolean functions requiring large circuits can be shown by a non-constructive counting argument, as established by Shannon in 1949 [Sha49]. It follows from Shannon’s seminal result and a simple padding argument that if  $a > b \geq 1$  there are functions computable by circuits of size  $n^a$  that cannot be computed by circuits of size  $n^b$ . In other words, the classification of Boolean functions by their minimum circuit size forms a strict *hierarchy*.

Obtaining a “constructive” form of these results has been a holy grail in computational complexity theory for several decades due to its connections to derandomization and as an approach to separating P and NP. For instance, if there is a polynomial-time algorithm that given  $1^n$  outputs the truth-table of a function  $f: \{0, 1\}^{\log n} \rightarrow \{0, 1\}$  that requires circuits of size  $n^{\Omega(1)}$ , then  $P = BPP$  [IW97]. In results of this form, a constructive form of the (non-constructive) proof of the existence of hard functions is interpreted *computationally* as the existence of an algorithm of bounded complexity that computes a hard function.

In this paper, rather than focusing on the existence of algorithms to capture the constructiveness of a statement, we explore this notion from the perspective of mathematical logic, specifically concerning its *provability* in certain mathematical theories. We are interested in identifying the weakest theory capable of establishing the aforementioned circuit size hierarchy for Boolean circuits and related results.

As one of our contributions, we present a tight connection between the computational and proof-theoretic perspectives. We demonstrate that proving the non-uniform circuit size hierarchy in a theory known as  $T_2^1$  implies the existence of a function in  $P^{NP}$  that requires Boolean circuits of size at least  $n^{1+\epsilon}$ . The latter is a frontier question in complexity theory (see, e.g., [CMMW19]). Thus, in a precise sense, developing more constructive proofs of the circuit size hierarchy would lead to significant progress on explicit circuit lower bounds.

We now proceed to describe this result and other contributions of this work in detail.

## 1.2 Results

We will be concerned with standard theories of bounded arithmetic. These theories are designed to capture proofs that manipulate and reason with concepts from a specified complexity class. Notable examples include Cook’s theory  $PV_1$  [Coo75], which formalizes polynomial-time reasoning; Jeřábek’s theory  $APC_1$  [Jeř04, Jeř05, Jeř07], which extends  $PV_1$  by incorporating the dual weak pigeonhole principle for polynomial-time functions and formalizes probabilistic polynomial-time reasoning; and Buss’s theories  $T_2^i$  [Bus86], which incorporate induction principles corresponding to various levels of the polynomial-time hierarchy.

For an introduction to bounded arithmetic, we refer to [Bus97]. For its connections to computational complexity and a discussion on the formalization of complexity theory, we refer to [Oli24].<sup>1</sup> Here we only recall that theory  $PV_1$  corresponds essentially to  $T_2^0$  [Jeř06], and that  $T_2^0 \subseteq T_2^1 \subseteq T_2^2$  correspond to the first levels of Buss’s hierarchy. A brief overview of the theories is provided in Section 2.

For a given  $n \in \mathbb{N}$ , we use  $CIRCUIT[s(n)]$  to denote the set of Boolean functions  $f: \{0, 1\}^n \rightarrow \{0, 1\}$  computed by circuits of size at most  $s(n)$ . Similarly, we write  $FORMULA[s(n)]$  when referring to formula size. We use  $SIZE[s(n)]$  to denote the set of languages  $L \subseteq \{0, 1\}^*$  that admit a sequence of circuits of size at most  $s(n)$ .

---

<sup>1</sup>In particular, the reference [Oli24] contains a detailed discussion of some aspects of the formalization of the statements appearing below.

83 **Circuit Size Hierarchy.** For rationals  $a > b \geq 1$  and  $n_0$ , we consider the following sentence:<sup>2</sup>

$$\begin{aligned} \text{CSH}[a, b, n_0] \equiv & \forall n \geq n_0 \in \text{Log}, \exists \text{ circuit } D: \{0, 1\}^n \rightarrow \{0, 1\} \text{ of size } \leq n^a, \\ & \forall \text{ circuit } C: \{0, 1\}^n \rightarrow \{0, 1\} \text{ of size } \leq n^b, \exists x \in \{0, 1\}^n \text{ such that } D(x) \neq C(x). \end{aligned}$$

84 In other words,  $\text{CSH}[a, b, n_0]$  states that  $\text{CIRCUIT}[n^b] \subsetneq \text{CIRCUIT}[n^a]$  whenever  $n \geq n_0$ .

85

86 Next, we state our first result.

87 **Theorem 1.** *The following results hold:*

(i) For every choice of rationals  $a$  and  $b$  with  $a - 1 > b > 1$ , and for every large enough  $n_0 \in \mathbb{N}$ ,

$$\text{T}_2^2 \vdash \text{CSH}[a, b, n_0].$$

(ii) If there are rationals  $a > b > 1$  and a constant  $n_0 \in \mathbb{N}$  such that

$$\text{T}_2^1 \vdash \text{CSH}[a, b, n_0],$$

88 then there is a constant  $\varepsilon > 0$  and a language  $L \in \text{P}^{\text{NP}}$  such that  $L \notin \text{SIZE}[n^{1+\varepsilon}]$ .

89 (iii) Similarly to the previous item, if  $\text{PV}_1 \vdash \text{CSH}[a, b, n_0]$ , there is  $L \in \text{P}$  such that  $L \notin \text{SIZE}[n^{1+\varepsilon}]$ .

90 To put it another way, we can establish a circuit size hierarchy within the theory  $\text{T}_2^2$ . If this result could  
 91 also be proven in the theory  $\text{T}_2^1$ , it would lead to a significant breakthrough in circuit lower bounds. Thus, by  
 92 enhancing the proof complexity upper bound for the provability of the circuit size hierarchy, we can achieve  
 93 new circuit lower bounds.

94 Note that in Theorem 1 Items (ii) and (iii) we obtain a lower bound against circuits of size  $n^{1+\varepsilon}$ , where  
 95 the constant  $\varepsilon > 0$  depends on the proof of  $\text{CSH}[a, b, n_0]$  in the corresponding theory. In other words, while  
 96 the sentence claims the existence of hardness against circuits of size  $n^b$ , we are only able to extract a weaker  
 97 lower bound for an explicit problem.

98 In our next result, we describe a setting where we can extract all the hardness from a proof of the  
 99 corresponding sentence.

100 **Succinct Circuit Size Hierarchy.** For rationals  $a > b \geq 1$  and  $n_0$ , we consider the following sentence:

$$\begin{aligned} \text{SCSH}[a, b, n_0] \equiv & \forall n \geq n_0 \in \text{Log}, \exists \text{ collection } \{(x^1, b^1), \dots, (x^\ell, b^\ell)\} \text{ of size } \ell \leq n^a \text{ with} \\ & |x^i| = n \wedge |b^i| = 1 \text{ for each } i \in [\ell] \text{ and } x^i \neq x^j \text{ for distinct } i, j \in [\ell], \\ & \forall \text{ circuit } C: \{0, 1\}^n \rightarrow \{0, 1\} \text{ of size } \leq n^b, \exists i \in [\ell] \text{ such that } C(x^i) \neq b^i. \end{aligned}$$

101 In other words,  $\text{SCSH}[a, b, n_0]$  states that for every  $n \geq n_0$  there is a collection of  $\ell \leq n^a$  labelled examples  
 102 such that every circuit of size at most  $n^b$  disagrees with at least one of its labels.

103

104 We obtain the following results on the proof complexity of the succinct circuit size hierarchy.

105 **Theorem 2.** *The following results hold:*

---

<sup>2</sup>The abbreviation  $n \in \text{Log}$  denotes that  $n$  is the length of a variable  $N$  (see, e.g., [Oli24] for more details).

(i) For every choice of rationals  $a > b > 1$  and for every large enough  $n_0 \in \mathbb{N}$ ,

$$\mathsf{T}_2^2 \vdash \text{SCSH}[a, b, n_0].$$

(ii) If there are rationals  $a > b > 1$  and a constant  $n_0 \in \mathbb{N}$  such that

$$\mathsf{T}_2^1 \vdash \text{SCSH}[a, b, n_0],$$

106 then there is a language  $L \in \mathsf{P}^{\text{NP}}$  such that  $L \notin \text{SIZE}[n^b]$ .

107 In our final result, we investigate the provability of size hierarchies for more restricted computational  
108 models in  $\mathsf{T}_2^1$  and weaker theories.

109 **Formula Size Hierarchy.** For rationals  $a > b \geq 1$  and  $n_0$ , we consider the following sentence:

$$\begin{aligned} \text{FSH}[a, b, n_0] \equiv & \forall n \geq n_0 \in \text{Log}, \exists \text{ formula } F: \{0, 1\}^n \rightarrow \{0, 1\} \text{ of size } \leq n^a, \\ & \forall \text{ formula } G: \{0, 1\}^n \rightarrow \{0, 1\} \text{ of size } \leq n^b, \exists x \in \{0, 1\}^n \text{ such that } F(x) \neq G(x). \end{aligned}$$

110 In other words,  $\text{FSH}(a, b, n_0)$  states that  $\text{FORMULA}[n^b] \subsetneq \text{FORMULA}[n^a]$  whenever  $n \geq n_0$ .

111

112 We establish that for some parameters a formula size hierarchy is provable already in  $\text{PV}_1$ .

**Theorem 3.** Consider rationals  $a > 2$  and  $b = 3/2$ , and let  $n_0$  be a large enough positive integer. Then

$$\text{PV}_1 \vdash \text{FSH}[a, b, n_0].$$

113 While many lower bounds can be proven in  $\text{APC}_1$  and stronger theories (see [MP20, Oli24, CLO24] and  
114 references therein), Theorem 3 provides an example of a non-trivial lower bound that can be established in  
115  $\text{PV}_1$ , which might be of independent interest.

### 116 1.3 Techniques

117 The proofs of Items (ii) and (iii) in Theorem 1 are inspired by arguments from [KO17, Kra21] that  
118 rely on a combination of a witnessing theorem with a term elimination strategy. Recall that the witnessing  
119 theorem allows us to extract computational information from a proof of the sentence in the theory. Roughly  
120 speaking, in our context this implies that the first existential quantifier in the sentence  $\text{CSH}[a, b, n_0]$ , which  
121 corresponds to a circuit computing a hard function, can be witnessed by a finite number of terms  $t_1, \dots, t_k$   
122 of the corresponding theory. In  $\text{PV}_1$ , a term yields a polynomial-time function, while in  $\mathsf{T}_2^1$  a term yields  
123 a polynomial-time function with access to an NP oracle. The main difficulty is that (1) for a given input  
124 length  $n$  it is not clear which term among  $t_1, \dots, t_k$  succeeds in constructing a hard function, and (2) for a  
125 term to succeed we must provide counter-examples to the candidate witnesses provided by previous terms.

126 As in previous papers, we assume that the conclusion of the theorem does not hold, and use this assump-  
127 tion to rule out the correctness of each term. This leads to a contradiction, meaning that the original sentence  
128 is not provable in the corresponding theory. Implementing this plan requires a careful argument, and we are  
129 currently only able to carry it out under a complexity inclusion in  $\text{SIZE}[n^{1+\varepsilon}]$  as opposed to  $\text{SIZE}[n^b]$ . The  
130 proof of the result is given in Section 3.1.

131 On the other hand, in the case of the succinct circuit size hierarchy, the argument for Item (ii) of Theo-  
132 rem 2 is simpler and allows us to start with the weaker assumption that  $\mathsf{P}^{\text{NP}} \subseteq \text{SIZE}[n^b]$ . Without getting

133 into the technical details, the main reason for not losing hardness in this result is that given a labelled list of  
 134 examples and access to an NP oracle, we can efficiently compute a minimum size circuit that agrees with  
 135 this list of inputs. Consequently, we can check if a candidate labelled list provided by a term is indeed hard,  
 136 or produce a counter-example when this is not the case. The same computation is not available in the case  
 137 of Theorem 1, since it is not clear how to efficiently compute with access to an NP oracle if a given circuit  
 138 admits a smaller equivalent circuit. The proof of Item (ii) of Theorem 2 appears in Section 3.2.

139 The proofs of Theorem 1 Item (i) and Theorem 2 Item (i) are given in Section 3.3. The formalization of  
 140 these hierarchies in  $T_2^2$  is easily done with access to the dual Weak Pigeonhole Principle for polynomial-time  
 141 functions, a principle which is known to be available in  $T_2^2$ . In more detail, CSH follows from SCSH in  $PV_1$ ,  
 142 while SCSH can be established in theory  $APC_1$ , which is contained in  $T_2^2$ .

143 Finally, in the proof of Theorem 3 we formalize in  $PV_1$  that the parity function on  $n$  bits can be computed  
 144 by formulas of size  $O(n^2)$  and require formulas of size  $\Omega(n^{3/2})$ . This yields in  $PV_1$  a proof of  $FSH[a, b, n_0]$   
 145 for any choice of parameters  $a > 2$ , large enough  $n_0$ , and  $b = 3/2$ . The upper bound on the complexity  
 146 of parity follows from a straightforward formalization of the correctness of the formula obtained via a  
 147 divide-and-conquer procedure. On the other hand, in order to show the formula lower bound we formalize  
 148 Subbotovskaya's argument [Sub61] based on the method of restrictions. To implement the proof in  $PV_1$ , we  
 149 directly define an efficient refuter that given a small formula outputs an input string where it fails to compute  
 150 the parity function. The correctness of the refuter is established by induction using an induction principle  
 151 available in the theory  $S_2^1$ . We then rely on a conservation result showing that the proof can also be done in  
 152  $PV_1$ . A detailed exposition of the argument appears in Section 4.

153 **Acknowledgements.** We thank Emil Jeřábek for a discussion about witnessing theorems in bounded arith-  
 154 metic. This work received support from the Royal Society University Research Fellowship URF\R1\191059;  
 155 the UKRI Frontier Research Guarantee Grant EP/Y007999/1; and the Centre for Discrete Mathematics and  
 156 its Applications (DIMAP) at the University of Warwick.

## 157 2 Preliminaries

### 158 2.1 Complexity Theory

159 We employ standard definitions from complexity theory, such as basic complexity classes, Boolean  
 160 circuits, and Boolean formulas (see, e.g., [AB09]).

161 Let  $\mathbb{N}$  represent the set of non-negative integers. For any  $a \in \mathbb{N}$ , let  $|a|$  denote the length of its binary  
 162 representation, defined as  $|a| \triangleq \lceil \log_2(a + 1) \rceil$ . For a constant  $k \geq 1$ , a function  $f: \mathbb{N}^k \rightarrow \mathbb{N}$  is said to  
 163 be computable in polynomial time if  $f(x_1, \dots, x_k)$  can be computed in time polynomial in  $|x_1|, \dots, |x_k|$ .  
 164 For convenience, we might write  $|\vec{x}| \triangleq |x_1|, \dots, |x_k|$ . The class FP denotes the set of polynomial-time  
 165 computable functions. Although the definition of polynomial time typically refers to a machine model,  
 166 FP can also be defined in a machine-independent manner as the closure of a set of base functions  $\mathcal{F}$  (not  
 167 described here) under *composition* and *limited recursion on notation*. A function  $f(\vec{x}, y)$  is defined from  
 168 functions  $g(\vec{x})$ ,  $h(\vec{x}, y, z)$ , and  $k(\vec{x}, y)$  by *limited recursion on notation* if

$$\begin{aligned} f(\vec{x}, 0) &= g(\vec{x}) \\ f(\vec{x}, y) &= h(\vec{x}, y, f(\vec{x}, \lfloor y/2 \rfloor)) \\ f(\vec{x}, y) &\leq k(\vec{x}, y) \end{aligned}$$

169 for every sequence  $(\vec{x}, y)$  of natural numbers. Cobham [Cob65] established that FP is the smallest class  
 170 of functions that contains the base functions  $\mathcal{F}$  and is closed under composition and limited recursion on  
 171 notation.

## 172 2.2 Bounded Arithmetic

### 173 2.2.1 Logical Theories

174 We recall the definitions of some standard theories of bounded arithmetic. For more details, the reader  
 175 can consult [Kra95, CN10, Kra19].

176 **Cook’s Theory  $PV_1$  [Coo75].** The first-order theory PV is designed to model the set  $\mathbb{N}$  of natural numbers  
 177 with the standard interpretations for constants and function symbols like  $0, +, \times$ , etc. The vocabulary (lan-  
 178 guage) of PV, denoted  $\mathcal{L}_{PV}$ , includes a function symbol for each polynomial-time algorithm  $f: \mathbb{N}^k \rightarrow \mathbb{N}$ ,  
 179 where  $k$  is any constant. These function symbols and their defining axioms are derived using Cobham’s char-  
 180 acterization of polynomial-time functions discussed above. PV also includes an induction axiom scheme  
 181 that simulates binary search, and it can be shown that it allows induction over quantifier-free formulas (i.e.,  
 182 polynomial-time predicates).

183 PV can be formulated with all axioms as universal formulas (i.e.,  $\forall \vec{x} \phi(\vec{x})$ , where  $\phi$  is free of quanti-  
 184 fiers). Thus, PV is a *universal theory*. Although the definition of PV is quite technical, the theory is fairly  
 185 robust and the details of its definition are often unnecessary for practical purposes. In particular, PV has an  
 186 equivalent formalization that does not rely on Cobham’s result [Jeř06].

**Jeřábek’s Theory  $APC_1$  [Jeř04, Jeř05, Jeř07].**  $APC_1$  extends PV with the *dual Weak Pigeonhole Prin-*  
*ciple* (dWPHP) for PV functions:

$$APC_1 \triangleq PV \cup \{dWPHP(f) \mid f \in \mathcal{L}(PV)\}.$$

187 Each sentence  $dWPHP(f)$  postulates that, for every length  $n = |N|$  and for every choice of  $\vec{z}$ , there is  
 188  $y < (1 + 1/n) \cdot 2^n$  such that  $f(\vec{z}, x) \neq y$  for every  $x < 2^n$ . It is known that  $APC_1$  is contained in  $T_2^2$   
 189 [MPW02].

190 **Buss’s Theories  $S_2^i$  and  $T_2^i$  [Bus86].** The language  $\mathcal{L}_B$  for these theories includes predicate symbols  $=$   
 191 and  $\leq$ , constant symbols 0 and 1, and function symbols  $S$  (successor),  $+$ ,  $\cdot$ ,  $\lfloor x/2 \rfloor$ ,  $|x|$  (interpreted as the  
 192 length of  $x$ ), and  $\#$  (interpreted as  $x \# y = 2^{|x| \cdot |y|}$ , known as “smash”).

193 Recall that a *bounded quantifier* is a quantifier of the form  $Qy \leq t$ , where  $Q \in \{\exists, \forall\}$  and  $t$  is a term  
 194 not involving  $y$ . Similarly, a *sharply bounded quantifier* is one of the form  $Qy \leq |t|$ . A formula where each  
 195 quantifier appears bounded (or sharply bounded) is called a bounded (or sharply bounded) formula.

196 We can create a hierarchy of formulas by counting alternations of bounded quantifiers. The class  $\Pi_0^b =$   
 197  $\Sigma_0^b$  contains the sharply bounded formulas. Recursively, for each  $i \geq 0$ , the classes  $\Sigma_i^b$  and  $\Pi_i^b$  are defined  
 198 by the quantifier structure of the sentence, ignoring sharply bounded quantifiers. For instance, if  $\varphi \in \Sigma_0^b$   
 199 and  $\psi \triangleq \exists y \leq t(\vec{x}) \varphi(y, \vec{x})$ , then  $\psi \in \Sigma_1^b$ . For the general case of the definition, see [Kra95]. It is known  
 200 that for each  $i \geq 1$ , a predicate  $P(\vec{x})$  is in  $\Sigma_i^p$  (the  $i$ -th level of the polynomial hierarchy) if and only if there  
 201 is a  $\Sigma_i^b$ -formula that agrees with it over  $\mathbb{N}$ .

202 These theories share a common set of finitely many axioms, BASIC, which postulate the expected  
 203 arithmetic behavior of the constants, predicates, and function symbols. The only difference among the  
 204 theories is the type of induction axiom scheme each one postulates.

$T_2^i$  is a theory in the language  $\mathcal{L}_B$  that extends BASIC by including the induction axiom IND:

$$\varphi(0) \wedge \forall x (\varphi(x) \rightarrow \varphi(x+1)) \rightarrow \forall x \varphi(x)$$

205 for all  $\Sigma_i^b$ -formulas  $\varphi(a)$ . The formula  $\varphi(a)$  may contain other free variables in addition to  $a$ .

$S_2^i$  is a theory in the language  $\mathcal{L}_B$  that extends BASIC by including the polynomial induction axiom PIND:

$$\varphi(0) \wedge \forall x (\varphi(\lfloor x/2 \rfloor) \rightarrow \varphi(x)) \rightarrow \forall x \varphi(x)$$

206 for all  $\Sigma_i^b$ -formulas  $\varphi(a)$ . The formula  $\varphi(a)$  may contain other free variables in addition to  $a$ .

207 **Theorem  $S_2^1$ (PV).** When proving some results in  $S_2^1$ , it is often convenient to use a more expressive vo-  
 208 cabulary that easily describes any polynomial-time function. This can be done in a *conservative* manner,  
 209 meaning the power of the theory is not increased. Specifically, let  $\Gamma$  be a set of  $\mathcal{L}_B$ -formulas. We say that  
 210 a polynomial-time function  $f: \mathbb{N}^k \rightarrow \mathbb{N}$  is  $\Gamma$ -*definable* in  $S_2^1$  if there exists a formula  $\psi(\vec{x}, y) \in \Gamma$  such that  
 211 the following conditions are met:

212 (i) For every  $\vec{a} \in \mathbb{N}^k$ ,  $f(\vec{a}) = b$  if and only if  $\mathbb{N} \models \varphi(\vec{a}, b)$ .

213 (ii)  $S_2^1 \vdash \forall \vec{x} (\exists y (\varphi(\vec{x}, y) \wedge \forall z (\varphi(\vec{x}, z) \rightarrow y = z)))$ .

214 Every function  $f \in \text{FP}$  is  $\Sigma_1^b$ -definable in  $S_2^1$ . By incorporating all functions in FP into the vocabulary  
 215 of  $S_2^1$  and extending the axioms of  $S_2^1$  with their defining equations, we obtain a theory  $S_2^1(\text{PV})$ . This  
 216 theory allows polynomial-time predicates to be referred to using quantifier-free formulas.  $S_2^1(\text{PV})$  remains  
 217 conservative over  $S_2^1$ , meaning any  $\mathcal{L}_B$ -sentence provable in  $S_2^1(\text{PV})$  is also provable in  $S_2^1$ . Finally, it is  
 218 known that  $S_2^1(\text{PV})$  proves the polynomial induction scheme for both  $\Sigma_1^b$ -formulas and  $\Pi_1^b$ -formulas within  
 219 the extended vocabulary.

## 220 2.2.2 The KPT Witnessing Theorem

221 The following witnessing theorem (a variant of Herbrand's theorem) is proved in [KPT91] (cf. also  
 222 [Kra95, Theorem 7.4.1]) for universal theories (like the theory  $\text{PV}_1$ ).

**Theorem 4** (KPT Theorem for  $\forall\exists\forall\exists$  sentences). *Let  $T$  be a universal theory with vocabulary  $\mathcal{L}$ . Let  $\varphi$  be an open  $\mathcal{L}$ -formula, and suppose that*

$$T \vdash \forall x \exists y \forall z \exists w \varphi(x, y, z, w).$$

*Then there is a finite sequence  $s_1, \dots, s_k$  of  $\mathcal{L}$ -terms such that*

$$T \vdash \forall x, z_1, \dots, z_k (\psi(x, s_1(x), z_1) \vee \psi(x, s_2(x, z_1), z_2) \vee \dots \vee \psi(x, s_k(z_1, \dots, z_{k-1}), z_k)),$$

*where*

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

223 For completeness, we describe a proof of Theorem 4 in Appendix A.

224 We can also apply the KPT Theorem to each theory  $T_2^i$  (for  $i \geq 1$ ) using a conservative extension of  
 225 the theory that admits a universal axiomatization. The corresponding theory is called  $\text{PV}_{i+1}$  [KPT91]. In  
 226  $\text{PV}_{i+1}$ , each term is equivalent to an  $\text{FP}^{\Sigma_i^p}$  function over the standard model. This leads to the following  
 227 result.



**Theorem 5** (Consequence of the KPT Theorem for Theory  $T_2^i$ ). *Let  $i \geq 1$ ,  $\varphi(x, y, w, z)$  be a  $\Pi_i^b$ -formula, and suppose that*

$$T_2^i \vdash \forall x \exists y \forall z \exists w \varphi(x, y, w, z).$$

*Then there is a finite sequence  $f_1, \dots, f_k$  of function symbols, each corresponding to an  $FP^{\Sigma_i^p}$  function, such that*

$$\mathbb{N} \models \forall x, z_1, \dots, z_k (\psi(x, f_1(x), z_1) \vee \psi(x, f_2(x, z_1), z_2) \vee \dots \vee \psi(x, f_k(z_1, \dots, z_{k-1}), z_k)),$$

where

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

## 228 3 Circuit Size Hierarchies in Bounded Arithmetic

### 229 3.1 Explicit Circuit Lower Bounds from Provability in $PV_1$ and $T_2^1$

230 In this section, we prove Theorem 1 Items (ii) and Items (iii).

**Theorem 6** (Theorem 1 Item (iii)). *If there are rationals  $a > b > 1$  and  $n_0 \in \mathbb{N}$  such that*

$$PV_1 \vdash CSH[a, b, n_0],$$

231 *then there is a constant  $\varepsilon > 0$  and a language  $L \in P$  such that  $L \notin \text{SIZE}[n^{1+\varepsilon}]$ .*

232 *Proof.* Towards a contradiction, suppose that  $PV_1 \vdash CSH[a, b, n_0]$  for rationals  $a > b > 1$  and some  
233 constant  $n_0$  and that  $P \subseteq \bigcap_{\varepsilon > 0} \text{SIZE}[n^{1+\varepsilon}]$ . The sentence  $CSH[a, b, n_0]$  has the form  $\forall \exists \forall \exists$ :

$$CSH[a, b, n_0] \triangleq \forall n \geq n_0 \in \text{Log}, \exists \text{circuit } D \forall \text{circuit } C \psi_{a,b}(n, D, C),$$

234 where  $\psi_{a,b}(n, D, C)$  is the existential formula:

$$\psi_{a,b}(n, D, C) \triangleq \exists x |x| \leq n \wedge \text{SIZE}(D) \leq n^a \wedge (\text{SIZE}(C) \leq n^b \rightarrow D(x) \neq C(x)).$$

235 Therefore, we can apply the KPT Theorem (Theorem 4), which provides  $PV_1$ -terms, equivalently FP func-  
236 tions,  $s_1, \dots, s_k$ , where  $k$  is a constant, such that

$$\mathbb{N} \models \psi_{a,b}(n, s_1(1^{(n)}), C_1) \vee \psi_{a,b}(n, s_2(1^{(n)}), C_1), C_2) \vee \dots \vee \psi_{a,b}(n, s_k(1^{(n)}), C_1, \dots, C_{k-1}), C_k). \quad (1)$$

237 In the relation above the circuits  $C_1, \dots, C_k$  are universally quantified.

238 Next, we use  $P \subseteq \bigcap_{\varepsilon > 0} \text{SIZE}[n^{1+\varepsilon}]$  to refute each of these disjuncts. We start by considering the fol-  
239 lowing language,  $D$ -Eval:

240

**Input:** A string  $x$  and a sequence  $\langle C_1, C_2, \dots, C_r \rangle$  of  $r \leq k - 1$  circuits

1 Define  $n \triangleq |x|$ ;

2 Simulate  $s_{r+1}(1^{(n)}, C_1, \dots, C_r)$  and interpret the output as a Boolean circuit  $D: \{0, 1\}^n \rightarrow \{0, 1\}$ ;

241 // We assume w.l.o.g. that  $D$  is a valid  $n$ -bit circuit of size  
 $\leq n^a$ , since otherwise the disjunct is trivially false.

3 Evaluate  $D$  on input  $x$  and output the result.

**Algorithm 1:** The pseudocode of an algorithm that decides the language  $D$ -Eval.

$D$ -Eval is in P due to the fact that  $s_1, \dots, s_k \in \text{FP}$  and circuit evaluation is in FP. By our assumption on the circuit complexity of the complexity class P, for every input length  $m$  and every  $\varepsilon > 0$ ,  $D\text{-Eval} \in \text{SIZE}[m^{1+\varepsilon}]$ , so we can choose

$$\varepsilon_0 \triangleq b^{1/(2k)} - 1 > 0$$

and have  $D\text{-Eval} \in \text{SIZE}[m^{b^{1/(2k)}}]$ . We also define the constants

$$\epsilon_i \triangleq b^{i/k} \quad \text{and} \quad \delta_i \triangleq b^{(2i-1)/(2k)}$$

242 for  $i = 1, \dots, k$ . Note that  $\epsilon_i = (1 + \varepsilon_0)\delta_i$  and  $\delta_{i+1} > \epsilon_i$ .

243 We start by refuting  $\psi_{a,b}(n, s_1(1^{(n)}), C_1)$ . We consider inputs of the form  $x, \lambda$  to  $D\text{-Eval}$ , where  $\lambda$  is  
 244 the empty sequence. Then the input has length  $n + c$ , where  $c = O(\log n)$  accounts for the overhead in  
 245 the encoding of the input. We consider the circuit  $C_1^* \in \text{SIZE}[(n + c)^{1+\varepsilon_0}]$ , which evaluates as  $D\text{-Eval}$   
 246 on inputs of length  $n + c$ , and we fix the input variables not related to  $x$  to represent the empty sequence.  
 247 The resulting circuit has as input an  $n$ -bit string  $x$  and computes according to  $s_1(1^{(n)})$  by definition of the  
 248  $D\text{-Eval}$  algorithm. For sufficiently large  $n$ , we have that  $n + c \leq n^{\delta_1} \Rightarrow (n + c)^{1+\varepsilon_0} \leq n^{(1+\varepsilon_0)\delta_1} = n^{\epsilon_1}$ ,  
 249 therefore we have the circuit  $C_1^* \in \text{SIZE}[n^{\epsilon_1}]$  which agrees with the circuit  $s_1(1^{(n)})$  on all  $n$ -bit inputs.  
 250 Since  $\epsilon_1 \leq b$ , we have that  $\mathbb{N} \not\models \psi_{a,b}(n, s_1(1^{(n)}), C_1^*)$ .

251 We can apply a similar argument to the next disjunct using the aforementioned circuit  $C_1^*$ . In more  
 252 detail, we consider the input  $(x, \langle C_1^* \rangle)$  on  $D\text{-Eval}$ , which has length  $m = n + 9n^{\epsilon_1} \log(n^{\epsilon_1}) + c \leq n^{\delta_2}$  for  
 253 sufficiently large  $n$  due to  $\delta_2 > \epsilon_1$ , and a corresponding circuit  $C_2^* \in \text{SIZE}[m^{1+\varepsilon_0}]$  provided by the circuit  
 254 upper bound hypothesis. Similarly, we can fix the  $n^{\epsilon_1} \log(n^{\epsilon_1}) + c$  variables not related to the input string  
 255  $x$ . This provides an  $n$ -bit circuit  $C_2^* \in \text{SIZE}[n^{\epsilon_2}]$  that computes according to the circuit  $s_2(1^{(n)}, C_1^*)$ , due to  
 256 the definition of the  $D\text{-Eval}$  algorithm. Since  $\epsilon_2 < b$ , we have that  $\mathbb{N} \not\models \psi_{a,b}(n, s_2(1^{(n)}, C_1^*), C_2^*)$ .

257 Inductively, if we have circuits  $C_1^*, C_2^*, \dots, C_i^*$  for some  $i \leq k - 1$  of sizes at most  $n^{\epsilon_1}, n^{\epsilon_2}, \dots, n^{\epsilon_i}$ ,  
 258 respectively, we consider the input  $(x, \langle C_1^*, \dots, C_i^* \rangle)$  to  $D\text{-Eval}$ , which has length  $m = n + 9n^{\epsilon_1} \log(n^{\epsilon_1}) +$   
 259  $\dots + 9n^{\epsilon_i} \log(n^{\epsilon_i}) + c \leq n^{\delta_{i+1}}$  for sufficiently large  $n$ . Therefore, by taking a corresponding  $m^{1+\varepsilon_0}$ -  
 260 size circuit for  $D\text{-Eval}$  and fixing all the inputs except for  $x$ , we get the circuit  $C_{i+1}^* \in \text{SIZE}[n^{\epsilon_{i+1}}] \subseteq$   
 261  $\text{SIZE}[n^b]$  which agrees with the circuit  $s_{i+1}(1^{(n)}, C_1^*, \dots, C_i^*)$  on all  $n$ -bit inputs. Consequently,  $\mathbb{N} \not\models$   
 262  $\psi_{a,b}(n, s_{i+1}(1^{(n)}, C_1^*, \dots, C_i^*), C_{i+1}^*)$ .

263 Overall, we can refute all disjuncts in Equation (1), which gives us a contradiction. This completes the  
 264 proof.  $\square$

**Theorem 7** (Theorem 1 Item (ii)). *If there are rationals  $a > b > 1$  and  $n_0 \in \mathbb{N}$  such that*

$$\text{T}_2^1 \vdash \text{CSH}[a, b, n_0],$$

265 *then there is a constant  $\varepsilon > 0$  and a language  $L \in \text{P}^{\text{NP}}$  such that  $L \notin \text{SIZE}[n^{1+\varepsilon}]$ .*

266 *Proof.* In this case, provability in  $\text{T}_2^1$  provides by the KPT Theorem (Theorem 5) functions  $s_1, \dots, s_k$  which  
 267 are in  $\text{FP}^{\text{NP}}$  instead of FP as in the previous proof. Therefore, the algorithm  $D\text{-Eval}$  is in  $\text{P}^{\text{NP}}$  and we use  
 268 the upper bound  $\text{P}^{\text{NP}} \subseteq \bigcap_{\varepsilon > 0} \text{SIZE}[n^{1+\varepsilon}]$  to get a contradiction in the same way as above.  $\square$

269 Note that in the arguments above we have no control over the constant  $\varepsilon > 0$ . It depends on the  
 270 number of disjuncts obtained from the KPT Theorem, which depends on the supposed proof of the hierarchy  
 271 sentence.

## 272 3.2 Extracting All the Hardness from Proofs of a Succinct Hierarchy Theorem

273 In this section, we prove Theorem 2 Item (ii).

**Theorem 8** (Theorem 2 Item (ii)). *If there are rationals  $a > b > 1$  and a constant  $n_0 \in \mathbb{N}$  such that*

$$\mathbb{T}_2^1 \vdash \text{SCSH}[a, b, n_0],$$

274 *then there is a language  $L \in \text{P}^{\text{NP}}$  such that  $L \notin \text{SIZE}[n^b]$ .*

275 *Proof.* The main idea here is to use the proof of SCSH in order to define a Turing machine  $M$  which runs  
276 in polynomial time using an NP oracle and its language is hard against  $n^b$ -size circuits.

277 Starting from  $\mathbb{T}_2^1 \vdash \text{SCSH}[a, b, n_0]$ , we see that the structure of the sentence is  $\forall \exists \forall \exists$ :

$$\text{SCSH}[a, b, n_0] \triangleq \forall n \geq n_0 \in \text{Log}, \exists \text{ collection } \mathcal{F}, \forall \text{ circuit } C \phi_{a,b}(n, \mathcal{F}, C),$$

278 where  $\phi_{a,b}(n, \mathcal{F}, C)$  is the formula that states that  $\mathcal{F}$  is a collection  $\{(x^1, b^1), \dots, (x^\ell, b^\ell)\}$  with  $\ell \leq n^a$ ,  
279 where  $|x^i| = n$  and  $|b^i| = 1$ , and that if  $C$  is a circuit on  $n$  variables and of size  $\leq n^b$ , then there is some  
280  $i \in [\ell]$  such that  $C(x^i) \neq b^i$  (we can move the existential quantifier at the front of the formula).

281 Thus, by the KPT Theorem (Theorem 5), there are  $\text{FP}^{\text{NP}}$  functions  $f_1, \dots, f_k$ , where  $k$  is a fixed con-  
282 stant, such that

$$\mathbb{N} \models \phi_{a,b}(n, f_1(1^{(n)}), C_1) \vee \phi_{a,b}(n, f_2(1^{(n)}), C_1, C_2) \vee \dots \vee \phi_{a,b}(n, f_k(1^{(n)}), C_1, \dots, C_{k-1}, C_k). \quad (2)$$

283 From the relation above, we can see that one of the functions  $f_1, \dots, f_k$  will output a collection that  
284 refutes every circuit of size  $\leq n^b$ . If it is not  $f_1$ , then there is a counterexample circuit  $C_1$ , which is used as  
285 extra input in  $f_2$  and so on. Since  $f_1, \dots, f_k$  are in  $\text{FP}^{\text{NP}}$ , we can simulate this procedure in a  $\text{P}^{\text{NP}}$  Turing  
286 machine  $M$ :

287

**Input:** A bit-string  $x$

- 1 Define  $n \triangleq |x|$ ;
- 2 **for**  $i = 1, \dots, k$  **do**
- 3     Simulate  $f_i$  with input  $1^{(n)}$  and, if  $i > 1$ ,  $C_1, \dots, C_{i-1}$ . Interpret the output as a collection  
 $\mathcal{F} = \{(x^1, b^1), \dots, (x^\ell, b^\ell)\}$  with  $\ell = n^a$ ;
- 4     Check with an NP oracle whether there exists a circuit  $C$  of size  $\leq n^b$ , such that  $C(x^i) = b^i$  for  
 all  $i \in [\ell]$ ;
- 5     If not or if  $i = k$ , exit the for-loop with the current  $\mathcal{F}$ ;
- 6     If there is such a circuit, then use the NP oracle to find it and name it  $C_i$ .
- 7 **end**
- 8 If the pair  $(x, 1)$  is in the collection  $\mathcal{F}$ , then **accept**. Else **reject**.

288

**Algorithm 2:** The Turing machine  $M_{a,b}$ , whose language is hard for  $n^b$ -size circuits.

289 It is easy to see that the language  $L(M_{a,b})$  recognised by the Turing machine  $M_{a,b}$ , is in  $\text{P}^{\text{NP}}$ . It suffices  
290 to show that  $L(M_{a,b}) \notin \text{SIZE}[n^b]$ .

291 Consider a circuit  $C \in \text{SIZE}[n^b]$ . Also, assume that the for-loop in Algorithm 2 ends in the  $r$ -th it-  
292 eration with  $r \leq k$ . We fix the circuits  $C_1, C_2, \dots, C_{r-1}$  found by the algorithm. Then the formula  
293  $\phi_{a,b}(n, f_r(1^{(n)}), C_1, \dots, C_{r-1}, C)$  always holds. If  $r < k$  and  $C$  did not satisfy it, then the NP oracle  
294 would find  $C$  as a counterexample and it would continue to the  $(r + 1)$ -th iteration. If  $r = k$ , then by the

295 construction of  $C_1, C_2, \dots, C_{k-1}$ , the formulas  $\phi_{a,b}(n, f_i(1^{(n)}, C_1, \dots, C_{i-1}), C_i)$  for  $i < k$  do not hold,  
 296 which means by Equation (2) that  $\phi_{a,b}(n, f_k(1^{(n)}, C_1, \dots, C_{k-1}), C)$  is true.

297 Since  $\mathcal{F} \equiv f_r(1^{(n)}, C_1, \dots, C_{r-1})$ , from  $\phi_{a,b}(n, \mathcal{F}, C)$ , we get that there is some  $i \in [\ell]$ , such that  
 298  $C(x^i) \neq b^i$ . However, if  $b^i = 1$ , then  $x^i \in L(M_{a,b})$ , and if  $b^i = 0$ , then  $x^i \notin L(M_{a,b})$ . In both cases, the  
 299 circuit  $C$  fails to recognise the language  $L(M_{a,b})$ , and the proof is complete.  $\square$

### 300 3.3 Formalization in $T_2^2$

301 In this section, we prove Theorem 1 Item (i) and Theorem 2 Item (i). To achieve this, we show that  
 302 the succinct circuit size hierarchy is provable in  $APC_1$ , which is contained in  $T_2^2$ . We then observe that the  
 303 circuit size hierarchy is easily provable from the succinct circuit size hierarchy.

**Theorem 9.** *For every choice of rationals  $a > b > 1$  and for every large enough  $n_0 \in \mathbb{N}$ ,*

$$APC_1 \vdash SCSH[a, b, n_0].$$

304 *In particular,  $SCSH[a, b, n_0]$  is provable in  $T_2^2$ .*

305 *Proof.* We define the polynomial-time function,  $f$ , which takes as input the description of a circuit,  $C$ , of  
 306 size  $n^b$ , which means that the length of the description of  $C$  is  $9n^b \log n^b$ , and outputs a bit string  $y$  of length  
 307  $n^a$  with the property that for all  $i = 0, 1, \dots, n^a - 1$ ,  $y_i = C(i)$ .

308 The correctness of the polynomial-time algorithm  $f$  is provable in  $PV_1$ . In other words,

$$PV_1 \vdash \forall n \in \text{Log} (|x| \leq 9n^b \log n^b \wedge |y| \leq n^a) \rightarrow [ |f(x)| \leq n^a \wedge (f(x) = y \leftrightarrow \forall i < n^a y_i = \text{Eval}(x, i)) ]. \quad (3)$$

309 The quantifier  $\forall i \leq n^a$  is sharply bounded, so this formula is provable in  $PV_1$ .

310 The theory  $APC_1$  includes the dWPHP axiom for all PV functions with input length  $n$  and output length  
 311  $n + 1$ , or equivalently input length  $n$  and output length  $m$  with  $n < m$ . From the first part of Equation (3),  
 312 the input length of  $f$  is  $9n^b \log n^b$ , while the output length is  $n^a$ . Furthermore, it is provable in  $PV_1$  that  
 313 there is some constant  $n_0$ , such that  $\forall n \geq n_0 n^a > 9n^b \log n^b$ . Therefore, we can use the axiom:

$$\text{dWPHP}(f) \triangleq \forall n \geq n_0 \exists y (|y| = n^a) \forall x (|x| = 9n^b \log n^b) f(x) \neq y \quad (4)$$

Every circuit of size  $n^b$  can be described by a string of size  $n^b \log n^b$ , which means that

$$\forall C \in \text{SIZE}[n^b] |C| \leq 9n^b \log n^b.$$

Also, from the second part of Equation (3), using the notation for the circuit  $C$ , we get that

$$f(C) \neq y \leftrightarrow \exists i < n^a C(i) \neq y_i.$$

314 Substituting the last two relations to Equation (4), we get that

$$APC_1 \vdash \forall n \geq n_0 \in \text{Log} \exists y (|y| = n^a) \forall C \in \text{SIZE}[n^b] \exists i < n^a C(i) \neq y_i,$$

315 which is equivalent with  $SCSH[a, b, n_0]$ .  $\square$

**Corollary 10.** *For every choice of rationals  $a > b > 1$  and for every large enough  $n_0 \in \mathbb{N}$ ,*

$$T_2^2 \vdash \text{CSH}[a + 1, b, n_0].$$

316 *Proof.* Since  $a > b$ , there is some rational  $\epsilon > 0$ , such that  $a - \epsilon > b$ . From Theorem 9, we know that  
 317  $\text{SCSH}[a - \epsilon, b, n_0]$  is provable in  $\text{APC}_1$  for every large enough  $n_0$ . Therefore, if we prove that

$$\text{PV}_1 \vdash \exists \text{ collection } \mathcal{F} = \{(x^1, b^1), \dots, (x^\ell, b^\ell)\} \text{ of size } \ell \leq n^{a-\epsilon} \text{ with } x^i \neq x^j \text{ for distinct } i, j \in [\ell] \rightarrow \\ \exists \text{ circuit } D: \{0, 1\}^n \rightarrow \{0, 1\} \text{ of size } \leq n^{a+1}, \forall i \in [\ell] D(x^i) = b^i,$$

318 we can easily deduce that  $\text{APC}_1 \vdash \text{CSH}[a + 1, b, n_0]$ . The same holds also for  $\text{T}_2^2$ .

Therefore, it is sufficient to argue in  $\text{PV}_1$  that there is a polynomial-time function  $\text{Circuit}(\mathcal{F})$ , which given the collection  $\mathcal{F}$ , outputs a circuit  $D: \{0, 1\}^n \rightarrow \{0, 1\}$  of the required size such that  $\forall i \in [\ell] D(x^i) = b^i$ . The construction of the circuit  $D$  is pretty straightforward: For every  $n$ -bit string  $x^i$ , such that  $(x^i, 1) \in \mathcal{F}$ , we construct the term  $T^i$ , which is the conjunction of the  $n$  bits of  $x^i$  (we put  $x_j$  if the  $j$ th bit of  $x^i$  is 1 and  $\neg x_j$  if the  $j$ th bit of  $x^i$  is 0). Then we make the DNF

$$D \triangleq \bigvee_{(x^i, 1) \in \mathcal{F}} T^i,$$

319 which agrees with the collection  $\mathcal{F}$ , and viewed as a circuit has size at most  $n^{a-\epsilon}(n+1)$  (at most  $n$   $\wedge$ -gates  
 320 for each one of the at most  $n^a$  terms and at most  $n^a$   $\vee$ -gates for the final disjunction). The correctness of the  
 321 resulting circuit is easily provable in  $\text{PV}_1$ , while for large enough  $n_0$ , we have  $\forall n \geq n_0 n^{a-\epsilon}(n+1) \leq n^{a+1}$ .  
 322 Hence, we have the desired result.  $\square$

### 323 3.4 On the Gap Between $\text{T}_2^1$ and $\text{T}_2^2$

324 We noticed above that it is possible to prove the circuit size hierarchy in the theory  $\text{T}_2^2$ . In contrast, it  
 325 seems difficult to implement a similar proof in the theory  $\text{T}_2^1$ . The reason behind this difficulty is connected  
 326 to the proof complexity of the dual Weak Pigeonhole Principle. If there is a proof of the circuit size hier-  
 327 archy in  $\text{T}_2^1$ , either it uses an approach that relies on a principle that is not equivalent to  $\text{dWPHP}(\text{PV})$ , or  
 328  $\text{dWPHP}(\text{PV})$  is also provable in  $\text{T}_2^1$ .

329 Paris, Wilkie, and Woods [PWW88] were the first to establish the provability of  $\text{dWPHP}(\text{PV})$  in Buss's  
 330 hierarchy. Subsequently, Maciel, Pitassi, and Woods [MPW02] provided an alternative proof with an explicit  
 331 inclusion of the principle in  $\text{T}_2^2$ . In this section, we explain why the same argument is not available in  $\text{T}_2^1$ .  
 332 (Their original proof is more general, and an exposition can be found in [Kra19].)

333 Assume that we have a PV-function  $g': \{0, 1\}^n \rightarrow \{0, 1\}^{n+1}$  with  $n \in \text{Log}$  or equivalently  $g': N \rightarrow$   
 334  $2N$ , such that  $\neg \text{dWPHP}_N^{2N}(g')$  holds. It is easy even in  $\text{S}_2^1(\text{PV})$  to extend this to a new function  $g: N \rightarrow N^2$ ,  
 335 such that  $\neg \text{dWPHP}_N^{N^2}(g) \triangleq \forall y < N^2 \exists x < N g(x) = y$  holds.

336 For  $\ell = 0, \dots, |N|$ , we consider all sequences  $w \in [N]^\ell$ . We extend a sequence by a new element  
 337 using the operation  $\frown$  (e.g.,  $(a_1, a_2, a_3) \frown a_4 = (a_1, a_2, a_3, a_4)$ ). For all sequences  $w$ , we define functions  
 338  $g_w: N/2^\ell \rightarrow N^2$  recursively as follows:

- 339 • If  $\ell = 0$ ,  $g_\emptyset = g$ .
- 340 • For  $i < N$ ,  $g_{w \frown i}(x) = y$  if  $\exists z < N$  such that  $g(z) = y \wedge g_w(x) = iN + z$ , otherwise output  $\emptyset$ . (Here  
 341  $\emptyset$  is just a fixed symbol that we use to denote “error” or that the function is undefined.)
- 342 •  $g_{w \frown N}(x) = y$  if  $\exists z < N \exists u < N$  such that  $g(z) = y \wedge g_w(x + N/2^{\ell+1}) = zN + u$ , otherwise output  
 343  $\emptyset$ .

344 Note that the formula  $g_w(x) = y$  is  $\Sigma_1^b$ -definable and that  $g_w(x)$  cannot have more than one value.

345 The key step of the proof is showing that

$$S_2^3 \vdash \neg \text{dWPHP}_N^{N^2}(g) \rightarrow \exists w \in [N]^\ell \neg \text{dWPHP}_{N/2^\ell}^{N^2}(g_w). \quad (5)$$

346 The right-hand size can be also written as

$$\exists w \in [N]^\ell \forall y < N^2 \exists x < N/2^\ell g_w(x) = y,$$

347 which is a  $\Sigma_3^b$  formula. Therefore, for the proof of Equation (5), we use  $\Sigma_3^b$ -LIND, which is available in  $S_2^3$ .

348 The intuition behind the inductive step is that if we split the domain into two equal intervals and the range  
349 into  $N$  intervals, from the surjectivity of  $g$  and  $g_w$ , either the first domain interval has all its values into the  
350  $i$ th range interval, which gives us the new sequence  $w \frown (i-1)$ , or the second domain interval has value at  
351 each one of the range intervals, which gives us the new sequence  $w \frown n$ .

352 To complete the argument, plugging  $\ell = |N|$  in Equation (5), we get a surjective function from 1 to  $N^2$ ,  
353 which is a clear contradiction when  $N > 1$ . Therefore,  $S_2^3 \vdash \text{dWPHP}(g)$ , and since  $S_2^3$  is  $\forall \Sigma_3^b$ -conservative  
354 over  $T_2^2$ , we also have  $T_2^2 \vdash \text{dWPHP}(g)$ .

355 The bottleneck to implement the proof in  $T_2^1$  is the quantifier complexity of the inductive statement  
356 associated with Equation (5). Another barrier for such a proof in  $T_2^1$  is the fact that for an arbitrary relation  
357  $R$ ,  $\text{dWPHP}(R)$  is not provable in  $S_2^2(R)$  [Kra92], so a proof of  $\text{dWPHP}(\text{PV})$  has to use some properties of  
358 PV functions.

## 359 4 Provability of Formula Size Bounds in $\text{PV}_1$

360 In this section, we prove Theorem 3. To achieve this, we establish that:

- 361 1. The parity function on  $n$  bits requires formulas of size  $\geq n^{3/2}$  (Section 4.1).
- 362 2. The parity function on  $n$  bits can be computed by formulas of size  $O(n^2) \leq n^a$  for any fixed rational  
363  $a > 2$  and large enough  $n$  (Section 4.2).
- 364 3. Consequently, the formula size hierarchy holds with parameters  $a > 2$  and  $b = 3/2$ , provided that  $n_0$   
365 is large enough (Section 4.3).

### 366 4.1 Subbotovskaya's Lower Bound

#### 367 4.1.1 High-Level Details of the Formalization

368 In this section, we sketch a formalization in  $\text{PV}_1$  of the proof that the parity function on  $n$  bits requires  
369 Boolean formulas of size  $\geq n^{3/2}$  [Sub61].<sup>3</sup> We adapt the argument presented in [Juk12, Section 6.3], which  
370 proceeds as follows:

1. [Juk12, Lemma 6.8]: Given a Boolean formula  $F$  on  $n$ -bit inputs, it is possible to fix one of its  
variables so that the resulting formula  $F_1$  satisfies

$$\text{Size}(F_1) \leq (1 - 1/n)^{3/2} \cdot \text{Size}(F).$$

---

<sup>3</sup>For concreteness, we let the size of a Boolean formula  $F$  be the number of leaves of  $F$  labeled by an input literal. We allow leaves that are labeled by constants, but we do not charge for them. Consequently, a constant function has formula complexity 0, while a non-constant function has formula complexity at least 1.

371 (In order to pick the variable to be restricted and its value, one first “normalizes” the formula  $F$ , as  
 372 implicitly described in [Juk12, Claim 6.9].)

2. [Juk12, Theorem 6.10]: By applying this result  $\ell \triangleq n - k$  times, it is possible to obtain a formula  $F_\ell$  on  $k$ -bit inputs such that

$$\text{Size}(F_\ell) \leq \text{Size}(F) \cdot (1 - 1/n)^{3/2} \cdot (1 - 1/(n-1))^{3/2} \dots (1 - 1/(k+1))^{3/2} = \text{Size}(F) \cdot (k/n)^{3/2}.$$

3. [Juk12, Example 6.11]: If the initial formula  $F$  computes the parity function, by setting  $\ell = n - 1$  we obtain

$$1 \leq \text{Size}(F_\ell) \leq (1/n)^{3/2} \cdot \text{Size}(F),$$

373 and consequently  $\text{Size}(F) \geq n^{3/2}$ .

374 We recommend reading this section with [Juk12, Section 6.3] at hand. We will slightly modify the  
 375 argument when formalizing the lower bound in  $\text{PV}_1$ . In more detail, given a small formula  $F$ , we recursively  
 376 construct (and establish correctness by induction) an  $n$ -bit input  $y$  witnessing that  $F$  does not compute  
 377 the parity function. (Actually, for technical reasons related to the induction step, we will simultaneously  
 378 construct an  $n$ -bit input  $y_n^0$  witnessing that  $F$  does not compute the parity function and an  $n$ -bit input  $y_n^1$   
 379 witnessing that  $F$  does not compute the negation of the parity function.)

380 Let  $s(n)$  be a size bound and  $\oplus(x)$  be a PV function that computes the parity of the binary string  
 381 described by  $x$ , i.e.,  $\oplus(x) \triangleq x_1 \oplus x_2 \oplus \dots \oplus x_n$ , where  $x_i$  denotes the  $i$ -th bit of  $x$ . To simplify notation,  
 382 we tacitly view  $x$  as a binary string. We assume that the formalization employs a well-behaved function  
 383 symbol  $\oplus$  such that  $\text{PV}_1$  proves the basic properties of the parity function, e.g.,  $\text{PV}_1 \vdash \oplus(x1) = 1 - \oplus(x)$   
 384 and  $\text{PV}_1 \vdash \oplus(x0) = \oplus(x)$ .

385 We consider the following  $\mathcal{L}(\text{PV})$ -sentence stating that the parity function requires formulas of size at  
 386 least  $s(n)$  for every input length  $n \geq 1$ :

$$\text{FLB}_s \triangleq \forall N \forall n \forall F (n = |N| \geq 1 \wedge \text{Size}(F) < s(n) \rightarrow \exists x (|x|_\ell = n \wedge \text{Eval}(F, x) \neq \oplus(x))),^4$$

387 where for convenience of notation we use the function symbol  $|w|_\ell$  to compute the bit-length of the string  
 388 represented by  $w$  (under some reasonable encoding).

389 **Theorem 11.** *Let  $s(n) \triangleq n^{3/2}$ . Then  $\text{PV}_1 \vdash \text{FLB}_s$ .*

390 *Proof.* Given  $b \in \{0, 1\}$ , we introduce the function  $\oplus^b(x) \triangleq \oplus(x) + b \pmod{2}$ . In order to prove  $\text{FLB}_s$  in  
 391  $\text{PV}_1$ , we explicitly consider a polynomial-time function  $R(1^n, F, b)$  with the following properties:<sup>5</sup>

- 392 1. Let  $b \in \{0, 1\}$ .  
 393 2. If  $\text{Size}(F) < s(n)$  then  $R(1^n, F, b)$  outputs an  $n$ -bit string  $y_n^b$  such that  $\text{Eval}(F, y_n^b) \neq \oplus^b(y_n^b)$ .

394 In other words,  $R(1^n, F, b)$  witnesses that the formula  $F$  does not compute the function  $\oplus^b$  over  $n$ -bit strings.  
 395 Note that the correctness of  $R$  is captured by the bounded universal sentence:

$$\text{Ref}_{R,s} \triangleq \forall 1^n \forall F (\text{Size}(F) < s(n) \rightarrow |y_n^0|_\ell = |y_n^1|_\ell = n \wedge F(y_n^0) \neq \oplus^0(y_n^0) \wedge F(y_n^1) \neq \oplus^1(y_n^1)),$$

<sup>4</sup>To simplify notation, we omit from the sentence  $\text{FLB}_s$  and in other parts of the exposition certain straightforward conditions, such as checking that  $F$  represents a valid formula and that it computes over  $n$ -bit input strings.

<sup>5</sup>For convenience, we often write  $1^n$  instead of explicitly considering parameters  $N$  and  $n = |N|$ . We might also write just  $F(x)$  instead of  $\text{Eval}(F, x)$ .

396 where we employed the abbreviations  $y_n^0 \triangleq R(1^n, F, 0)$  and  $y_n^1 \triangleq R(1^n, F, 1)$ . Our plan is to define  $R$  and  
 397 show that  $PV_1 \vdash \text{Ref}_{R,s}$ . Note that this implies  $\text{FLB}_s$  in  $PV_1$ . Jumping ahead, the correctness of  $R(1^n, F, b)$   
 398 will be established by polynomial induction on  $N$  (equivalently, induction on  $n = |N|$ ). Since  $\text{Ref}_{R,s}$  is a  
 399 universal sentence and  $S_2^1$  is  $\forall\Sigma_1^b$ -conservative over  $PV_1$ , polynomial induction for NP and coNP predicates  
 400 (admissible in  $S_2^1$ ; see, e.g., [Kra95, Section 5.2]) is available during the formalization. More details follow.

401 The procedure  $R(1^n, F, b)$  makes use of a few polynomial-time sub-routines (discussed below) and is  
 402 defined in the following way:

403

**Input:**  $1^n$  for some  $n \geq 1$ , formula  $F$  over  $n$ -bit inputs,  $b \in \{0, 1\}$ .

- 1 Let  $s(n) \triangleq n^{3/2}$ . If  $\text{Size}(F) \geq s(n)$  **return** “error”;
- 2 If  $\text{Size}(F) = 0$ ,  $F$  computes a constant function  $b_F \in \{0, 1\}$ . In this case, **return** the  $n$ -bit string  
 $y_n^b \triangleq y_1^b 0^{n-1}$  such that  $\oplus^b(y_1^b 0^{n-1}) \neq b_F$ ;
- 3 Let  $\tilde{F} \triangleq \text{Normalize}(1^n, F)$ ;  
 //  $\tilde{F}$  satisfies [Juk12, Claim 6.9],  $\text{Size}(\tilde{F}) \leq \text{Size}(F)$ ,  
 $\forall x \in \{0, 1\}^n F(x) = \tilde{F}(x)$ .
- 4 Let  $\rho \triangleq \text{Find-Restriction}(1^n, \tilde{F})$ , where  $\rho: [n] \rightarrow \{0, 1, \star\}$  and  $|\rho^{-1}(\star)| = n - 1$ ;  
 //  $\rho$  restricts a suitable variable  $x_i$  to a bit  $c_i$ , as in [Juk12,  
 Lemma 6.8].
- 5 Let  $F' \triangleq \text{Apply-Restriction}(1^n, \tilde{F}, \rho)$ . Moreover, let  $b' \triangleq b \oplus c_i$  and  $n' \triangleq n - 1$ ;  
 //  $F'$  is an  $n'$ -bit formula;  $\forall z \in \{0, 1\}^{\rho^{-1}(\star)} F'(z) = \tilde{F}(z \cup x_i \mapsto c_i)$ .
- 6 Let  $y_{n'}^{b'} \triangleq R(1^{n'}, F', b')$  and **return** the  $n$ -bit string  $y_n^b \triangleq y_{n'}^{b'} \cup y_i \mapsto c_i$ ;

404

**Algorithm 3:** Refuter Algorithm  $R(1^n, F, b)$ .

405 **Normalize( $1^n, F$ ) and its properties (in  $S_2^1$ ).** We say that a subformula  $G$  of  $F$  is a *neighbor* of a leaf  $z$   
 406 if either  $z \wedge G$  or  $z \vee G$  is a subformula of  $F$ . We say that a formula  $F$  over variables  $\{x_1, \dots, x_n\}$  is in  
 407 *normal form* if for every  $i \in [n]$  and every literal  $z \in \{x_i, \bar{x}_i\}$ , if  $z$  is a leaf of  $F$  and  $G$  is a neighbor of  $z$  in  
 408  $F$ , then  $G$  does not contain the variable  $x_i$ .

409 **Lemma 12.** *There is a polynomial-time function  $\text{Normalize}(1^n, F)$  that given a Boolean formula  $F$  over  $n$   
 410 input variables, outputs a formula  $\tilde{F}$  over  $n$  input variables such that the following holds:*

- 411 (i)  $\text{Size}(\tilde{F}) \leq \text{Size}(F)$ .
- 412 (ii) For every input  $x \in \{0, 1\}^n$ ,  $\tilde{F}(x) = F(x)$ .
- 413 (iii)  $\tilde{F}$  is in normal form.
- 414 (iv)  $\tilde{F}$  is either a constant 0 or 1, or  $\tilde{F}$  contains no leaves labeled by constants 0 and 1.

415 Moreover, the correctness of  $\text{Normalize}(1^n, F)$  is provable in  $S_2^1$ .

416 *Proof Sketch.* It is enough to verify that the proof of [Juk12, Claim 6.9] provides such a polynomial-time  
 417 function and that its correctness can be established in  $S_2^1$ . In more detail, if  $F$  is not in normal form, we can  
 418 efficiently compute a literal  $z \in \{x_i, \bar{x}_i\}$  and a neighbor  $G$  of  $z$  that violates the corresponding property.  
 419 As shown in [Juk12, Claim 6.9], we can fix any leaf  $z' \in \{x_i, \bar{x}_i\}$  in  $G$  by an appropriate constant  $c$  so  
 420 that the resulting formula  $F_1$  satisfies conditions (i) and (ii) of Lemma 12. After at most  $\ell \triangleq \text{Size}(F)$



421 iterations, we obtain a sequence  $F_1, \dots, F_\ell$  of formulas such that  $\tilde{F} \triangleq F_\ell$  satisfies conditions (i), (ii), and  
 422 (iii) of the lemma. Moreover, condition (iv) can always be guaranteed by simplifying the final formula,  
 423 i.e., by replacing subformulas  $0 \vee G$  by  $G$ ,  $1 \vee G$  by  $1$ ,  $0 \wedge G$  by  $0$ , and  $1 \wedge G$  by  $G$ . The correctness of  
 424  $\tilde{F} \triangleq \text{Normalize}(1^n, F)$  can be established by polynomial induction for coNP predicates (i.e.,  $\Pi_1^b$  formulas),  
 425 which is available in  $S_2^1$ .  $\square$

426 **Find-Restriction( $1^n, \tilde{F}$ ) and its properties (in  $S_2^1$ ).** We argue in  $S_2^1$  and follow the argument from the  
 427 proof of [Juk12, Lemma 6.8]. Let  $\tilde{F}$  be a formula over  $n$  input variables in normal form. We focus on the  
 428 non-trivial case, and assume that  $n \geq 2$ ,  $\text{Size}(\tilde{F}) \geq 2$ , and that  $\tilde{F}$  contains no leaves labeled by constants.  
 429 Let  $\text{Count}(1^n, F, i)$  be a polynomial-time algorithm that outputs the number of leaves of  $F$  that contain  
 430 the variable  $x_i$  (including its appearances as  $\bar{x}_i$ ). Let  $w = (w_1, \dots, w_n)$  be the corresponding sequence of  
 431 multiplicities, i.e.,  $w_i \triangleq \text{Count}(1^n, F, i)$ . Note that  $\sum_i w_i = \tilde{s}$ , where  $\tilde{s} \triangleq \text{Size}(\tilde{F})$ .

432 We claim that  $S_2^1$  proves the existence of an index  $i \in [n]$  such that  $w_i \geq \tilde{s}/n$ . First, for each  $j \in [n]$ , we  
 433 define the cumulative sum  $v_j \triangleq \sum_{i \leq j} w_i$ . Let  $v \triangleq (v_0, v_1, \dots, v_n)$  be the corresponding sequence, where  
 434 we set  $v_0 \triangleq 0$ . Notice that  $v_n = \tilde{s}$ . Since  $v$  contains  $n + 1$  elements, it can be efficiently computable from  
 435  $w$ . We now argue by induction on  $n$  that for some index  $j \in [n]$  we have  $v_j - v_{j-1} \geq v_n/n$ . This implies  
 436 that  $w_j = v_j - v_{j-1} \geq v_n/n = \tilde{s}/n$ , as desired.

437 If  $n = 1$ , then  $v_1 - v_0 = v_1 = v_1/1$  and the result holds for  $j = 1$ . Assume the result holds for  $n - 1$ , and  
 438 consider  $v_n$ . If  $v_n - v_{n-1} \geq v_n/n$ , we can pick  $j = n$  and we are done. Otherwise,  $v_{n-1} \geq v_n - v_n/n =$   
 439  $v_n(n - 1)/n$ . By the induction hypothesis, there is an index  $j \in [n - 1]$  such that  $v_j - v_{j-1} \geq v_{n-1}/(n - 1)$ .  
 440 Using the lower bound on  $v_{n-1}$ , we get that  $v_j - v_{j-1} \geq v_n/n$ , which concludes the proof.

Consequently,  $S_2^1$  proves the existence of a variable  $x_i$  which appears  $t \geq \tilde{s}/n$  times as a leaf of  $\tilde{F}$ . Let  
 $z_1, \dots, z_t$  be the leaves of  $\tilde{F}$  labeled by either  $x_i$  or  $\bar{x}_i$ . Recall that we assume that  $n \geq 2$ ,  $\text{Size}(\tilde{F}) \geq 2$ , and  
 that  $\tilde{F}$  satisfies conditions (iii) and (iv) of Lemma 12. Therefore, each leaf  $z_j$  has a neighbor subformula  $G_j$   
 in  $\tilde{F}$  that contains some leaf labeled by a literal not in  $\{x_i, \bar{x}_i\}$ . For this reason, if we set  $x_i$  to an appropriate  
 constant  $c_j$ ,  $G_j$  will disappear from  $F$ , thereby erasing at least another leaf not among  $z_1, \dots, z_t$ . As in  
 the proof of [Juk12, Lemma 6.8], if we let  $c \in \{0, 1\}$  be the constant that appears more often among  
 $c_1, \dots, c_t$  and set  $x_i \mapsto c$  in the restriction  $\rho$ , all the leaves  $z_1, \dots, z_t$  will be eliminated from  $\tilde{F}$  together  
 with at least  $t/2$  additional leaves.<sup>6</sup> Thus the total number of eliminated leaves, which we specify using a  
 polynomial-time function  $\text{NumRemoved}(1^n, \tilde{F}, \rho)$ , satisfies

$$\text{NumRemoved}(1^n, \tilde{F}, \rho) \geq t + \frac{t}{2} \geq \frac{3\tilde{s}}{2n}.$$

Overall, it follows that

$$S_2^1 \vdash \tilde{F} = \text{Normalize}(1^n, F) \wedge \rho = \text{Find-Restriction}(1^n, \tilde{F}) \rightarrow \text{NumRemoved}(1^n, \tilde{F}, \rho) \geq \frac{3}{2n} \cdot \text{Size}(\tilde{F}).$$

441 **Apply-Restriction( $1^n, \tilde{F}, \rho$ ) and its properties (in  $S_2^1$ ).** We only sketch the details. This is simply a  
 442 polynomial-time algorithm that, given a formula  $\tilde{F}$  on  $n$  input variables and a restriction  $\rho: [n] \rightarrow \{0, 1, *\}$   
 443 with  $|\rho^{-1}(\star)| = n - 1$  (i.e.,  $\rho$  restricts a single variable  $x_i$  to a constant  $c_i \in \{0, 1\}$ ), outputs a formula  $F'$   
 444 over  $n - 1$  input variables that sets every literal  $z \in \{x_i, \bar{x}_i\}$  to the corresponding constant and simplifies

<sup>6</sup>The existence of such a constant  $c$  can be proved in  $S_2^1$  in a way that is similar to the proof that some variable  $x_i$  appears in at least  $\tilde{s}/n$  leaves.

445 the resulting formula, e.g., replaces subformulas  $0 \vee G$  by  $G$ ,  $1 \vee G$  by  $1$ ,  $0 \wedge G$  by  $0$ , and  $1 \wedge G$  by  $G$ .  
 446 Additionally, for  $F' = \text{Apply-Restriction}(1^n, \tilde{F}, \rho)$ , we have

$$S_2^1 \vdash \text{Size}(F') \leq \text{Size}(\tilde{F}) - \text{NumRemoved}(1^n, \tilde{F}, \rho) \wedge \forall z \in \{0, 1\}^{\rho^{-1}(\star)} F'(z) = \tilde{F}(z \cup x_i \mapsto c_i). \quad (6)$$

447 Using the previously computed bound on  $\text{NumRemoved}(1^n, \tilde{F}, \rho)$  for  $\rho = \text{Find-Restriction}(1^n, \tilde{F})$ , we  
 448 obtain that for  $\tilde{F}$  and  $F'$  defined as above (with  $s' \triangleq \text{Size}(F')$  and  $\tilde{s} \triangleq \text{Size}(\tilde{F})$ ), and assuming that  $n \geq 2$ ,

$$S_2^1 \vdash s' \leq \tilde{s} - \frac{3}{2n} \cdot s' = \tilde{s} \cdot \left(1 - \frac{3}{2n}\right) \leq \tilde{s} \cdot \left(1 - \frac{1}{n}\right)^{3/2}. \quad (7)$$

449 The last inequality uses that  $S_2^1 \vdash \forall a, a \geq 2 \rightarrow (1 - 3/(2a))^2 \leq (1 - 1/a)^3$ , which one can easily verify.

450

451 Note that  $R(1^n, F, b)$  runs in time polynomial in  $n + |F| + |b|$  and that it is definable in  $S_2^1$ . Next, we  
 452 establish the correctness  $R(1^n, F, b)$  in  $S_2^1$ .

453 **Lemma 13.** *Let  $s(n) \triangleq n^{3/2}$ . Then  $S_2^1 \vdash \text{Ref}_{R,s}$ .*

*Proof.* We consider the formula  $\varphi(N)$  defined as

$$\forall F \forall n (n = |N| \wedge n \geq 1 \wedge \text{Size}(F) < s(n)) \rightarrow (|y_n^0|_\ell = |y_n^1|_\ell = n \wedge F(y_n^0) \neq \oplus^0(y_n^0) \wedge F(y_n^1) \neq \oplus^1(y_n^1)),$$

where as before we use  $y_n^0 \triangleq R(1^n, F, 0)$  and  $y_n^1 \triangleq R(1^n, F, 1)$ . Note that  $\varphi(N)$  is a  $\Pi_1^b$  formula. Below,  
 we argue that

$$S_2^1 \vdash \varphi(1) \quad \text{and} \quad S_2^1 \vdash \forall N \varphi(\lfloor N/2 \rfloor) \rightarrow \varphi(N).$$

454 Then, by polynomial induction for  $\Pi_1^b$  formulas (available in  $S_2^1$ ) and using that  $\varphi(0)$  trivially holds, it  
 455 follows that  $S_2^1 \vdash \forall N \varphi(N)$ . In turn, this yields  $S_2^1 \vdash \text{Ref}_{R,s}$ .

**Base Case:**  $S_2^1 \vdash \varphi(1)$ . In this case, for a given formula  $F$  and length  $n$ , the hypothesis of  $\varphi(1)$  is satisfied  
 only if  $n = 1$  and  $\text{Size}(F) = 0$ . Let  $y_1^0 \triangleq R(1, F, 0)$  and  $y_1^1 \triangleq R(1, F, 1)$ . We need to prove that

$$|y_1^0|_\ell = |y_1^1|_\ell = 1 \wedge F(y_1^0) \neq \oplus^0(y_1^0) \wedge F(y_1^1) \neq \oplus^1(y_1^1).$$

456 Since  $n = 1$  and  $\text{Size}(F) = 0$ ,  $F$  evaluates to a constant  $b_F$  on every input bit. The statement above is  
 457 implied by Line 2 in the definition of  $R(n, F, b)$ .

458 **(Polynomial) Induction Step:**  $S_2^1 \vdash \forall N \varphi(\lfloor N/2 \rfloor) \rightarrow \varphi(N)$ . Fix an arbitrary  $N$ , let  $n \triangleq |N|$ , and  
 459 assume that  $\varphi(\lfloor N/2 \rfloor)$  holds. By the induction hypothesis, for every formula  $F'$  with  $\text{Size}(F') < n^{3/2}$ ,  
 460 where  $n' \triangleq n - 1$ , we have

$$|y_{n'}^0|_\ell = |y_{n'}^1|_\ell = n' \wedge F'(y_{n'}^0) \neq \oplus^0(y_{n'}^0) \wedge F'(y_{n'}^1) \neq \oplus^1(y_{n'}^1), \quad (8)$$

461 where  $y_{n'}^0 \triangleq R(1^{n'}, F', 0)$  and  $y_{n'}^1 \triangleq R(1^{n'}, F', 1)$ .

Now let  $n \geq 2$ , and let  $F$  be a formula over  $n$ -bit inputs of size  $< n^{3/2}$ . By the size bound on  $F$ ,  
 $R(1^n, F, b)$  ignores Line 1. If  $\text{Size}(F) = 0$ , then similarly to the base case it is trivial to check that the con-  
 clusion of  $\varphi(N)$  holds. Therefore, we assume that  $\text{Size}(F) \geq 1$  and  $R(1^n, F, b)$  does not stop at Line 2. Let  
 $\tilde{F} \triangleq \text{Normalize}(1^n, F)$  (Line 3),  $\rho \triangleq \text{Find-Restriction}(1^n, \tilde{F})$  (Line 4),  $F' \triangleq \text{Apply-Restriction}(1^n, \tilde{F}, \rho)$   
 (Line 5),  $n' \triangleq n - 1$  (Line 5), and  $b' \triangleq b \oplus c_i$  (Line 5), where  $\rho$  restricts the variable  $x_i$  to the bit  $c_i$ .

Moreover, for convenience, let  $s \triangleq \text{Size}(F)$ ,  $\tilde{s} \triangleq \text{Size}(\tilde{F})$ , and  $s' \triangleq \text{Size}(F')$ . By Lemma 12 Item (i), Equation (7), and the bound  $s < n^{3/2}$ ,

$$S_2^1 \vdash s' \leq \tilde{s} \cdot (1 - 1/n)^{3/2} \leq s \cdot (1 - 1/n)^{3/2} < n^{3/2} \cdot (1 - 1/n)^{3/2} = (n - 1)^{3/2}.$$

462 Thus  $F'$  is a formula on  $n'$ -bit inputs of size  $< n'^{3/2}$ . Recall that for a given  $b \in \{0, 1\}$  we have  $b' = b \oplus c_i$ .  
 463 Let  $y_{n'}^{b'} \triangleq R(1^{n'}, F', b')$  (Line 6). By the first condition in the induction hypothesis (Equation (8)) and  
 464 the definition of each  $y_n^b \triangleq y_{n'}^{b'} \cup y_i \mapsto c_i$ , we have  $|y_n^0|_\ell = |y_n^1|_\ell = n$ . Below, we also rely on the last  
 465 two conditions in the induction hypothesis (Equation (8)), Lemma 12 Item (ii), and the last condition in  
 466 Equation (6). We derive the following statements, where  $b \in \{0, 1\}$ :

$$\begin{aligned} F'(y_{n'}^{b'}) &\neq \oplus^{b'}(y_{n'}^{b'}), \\ F(y_n^b) &= F'(y_{n'}^{b'}), \\ F(y_n^b) &\neq \oplus^{b'}(y_{n'}^{b'}). \end{aligned}$$

Notice that

$$\oplus^{b'}(y_{n'}^{b'}) = \oplus^{b \oplus c_i}(y_{n'}^{b'}) = c_i \oplus (\oplus^b(y_{n'}^{b'})) = c_i \oplus (\oplus^b(y_n^b) \oplus c_i) = \oplus^b(y_n^b).$$

467 These statements imply that, for each  $b \in \{0, 1\}$ ,  $F(y_n^b) \neq \oplus^b(y_n^b)$ . In other words, the conclusion of  $\varphi(N)$   
 468 holds. This completes the proof of the induction step.  $\square$

469 As explained above, the provability of  $\text{Ref}_{R,s}$  in  $S_2^1$  implies its provability in  $\text{PV}_1$ . Since  $\text{PV}_1 \vdash$   
 470  $\text{Ref}_{R,s} \rightarrow \text{FLB}_s$ , this completes the proof of Theorem 11.  $\square$

#### 471 4.1.2 On the Low-Level Details of the Formalization

472 In order to make our presentation accessible to a broader audience, in this section we provide more  
 473 details about the formalization of algorithms and about the proofs of their basic properties. For concreteness  
 474 and convenience, we consider the theory  $S_2^1(\text{PV})$ , i.e.,  $S_2^1$  extended with function symbols and axioms for  
 475 all polynomial-time functions as in Cobham's characterization of efficient computations. Since this theory  
 476 is  $\forall\Sigma_1^b$ -conservative over  $\text{PV}_1$  (see Section 2.2.1), the provability of  $\text{FLB}_s$  in  $S_2^1(\text{PV})$  yields its provability  
 477 in  $\text{PV}_1$ .

478 As a concrete example, we elaborate on a sub-routine employed by some algorithms discussed in Sec-  
 479 tion 4.1. We consider a polynomial-time function  $\text{Fix}(1^n, F, i, b)$  that, given the description of a formula  $F$   
 480 over  $n$  input variables, a variable index  $i \in [n]$ , and a bit  $b \in \{0, 1\}$ , replaces every leaf of  $F$  labeled by  $x_i$   
 481 with  $b$  and every leaf of  $F$  labeled by  $\bar{x}_i$  with  $1 - b$ , then returns the corresponding restricted formula  $F'$   
 482 over  $n - 1$  input variables (without the application of formula simplification rules). Next, we provide more  
 483 details about the specification of the procedure  $\text{Fix}$  in  $S_2^1(\text{PV})$  and about a proof of its correctness, i.e.,

$$S_2^1(\text{PV}) \vdash \forall 1^n \forall F \forall F' \forall x \forall z \forall i \tag{9}$$

$$(n \geq 2 \wedge |x|_\ell = n \wedge |z|_\ell = n - 1 \wedge 1 \leq i \leq n \wedge F' = \text{Fix}(1^n, F, i, b)) \rightarrow (\text{Eval}(F', z) = \text{Eval}(F, z \cup x_i \mapsto b)),$$

484 where  $z \cup x_i \mapsto b$  denotes a function that takes  $(z, i, b)$ , where  $z$  assigns bits to  $x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n$ ,  
 485 and outputs the  $n$ -bit string that agrees with  $z$  and sets  $x_i$  to  $b$ .

486 **On the Specification of  $\text{Fix}(1^n, F, i, b)$  in  $S_2^1(\text{PV})$ .** Theory  $S_2^1(\text{PV})$  contains function symbols for all  
487 polynomial-time algorithms according to Cobham's characterization of polynomial-time computations. Con-  
488 sequently, to specify  $\text{Fix}(1^n, F, i, b)$  we employ a definition of this computation in Cobham's formalism, i.e.,  
489 we define  $\text{Fix}(1^n, F, i, b)$  using simple base functions together with composition and recursion on notation.  
490 In order to be completely formal (a rather cumbersome task), one would first specify how formulas are rep-  
491 resented by numbers and the polynomial-time functions that manipulate the corresponding representation.  
492 We could then interpret the binary representation of an integer as two sequence of tuples, one describing  
493 the edges in the binary tree representation of the formula, and another describing the labels of each node  
494 of the tree. Finally,  $\text{Fix}(1^n, F, i, b)$  would be a routine that iterates over each leaf of  $F$  labelled by the  $i$ -th  
495 variable or its negation and replaces it with the appropriate constant. Using previously defined routines and  
496 their corresponding function symbols, a sequential algorithm of this form can be described as a recursive  
497 procedure in Cobham's characterization of polynomial-time functions. Moreover, we need to argue in the  
498 theory that the output length of the function on a given input is bounded by a polynomial, similarly to the  
499 constraint in the limited recursion on notation from Cobham's theorem.

500 **On the Proof of the Correctness of  $\text{Fix}(1^n, F, i, b)$  in  $S_2^1(\text{PV})$  (Equation (9)).**  $S_2^1(\text{PV})$  also contains  
501 axioms describing how the function symbols (polynomial-time functions) are obtained from each other.  
502 For instance,  $\text{Fix}(1^n, F, i, b)$  might use in its specification a routine  $R$  that takes as input a tuple describing a  
503 formula  $G$ , a bit  $b$ , and a leaf of  $G$  and its label, replaces the label of this leaf by the constant  $b$ , and outputs the  
504 new formula  $G'$ . We can then reason in  $S_2^1(\text{PV})$  about the correctness of  $\text{Fix}(1^n, F, i, b)$  (as in Equation (9))  
505 using the provable properties of  $R$  and of the function symbol  $\text{Eval}$ . In more detail,  $\text{Eval}$  can be defined  
506 recursively based on the structure of the input formula, and the base case of the proof of correctness relies  
507 on the properties of  $R$  and the fact that the internal evaluations of  $\text{Eval}(F', z)$  for  $F' = \text{Fix}(1^n, F, i, b)$  and  
508  $\text{Eval}(F, z \cup x_i \mapsto b)$  agree over all leaves. Crucially, the recursive nature of the specification of polynomial-  
509 time functions in Cobham's definition and in  $S_2^1(\text{PV})$  is compatible with the polynomial induction axioms  
510 available in  $S_2^1(\text{PV})$ , in the sense that we can define recursive procedures while simultaneously proving their  
511 relevant properties by induction.

## 512 4.2 Upper Bound

513 In this section, we show that the parity function on  $n$  bits can be computed by formulas of size  $O(n^2)$ ,  
514 provably in  $\text{PV}_1$ . We can formalize this upper bound in the language of  $\text{PV}$ , defining an  $\mathcal{L}(\text{PV})$ -sentence  
515 stating that the parity function can be computed by a formula of size  $s(n)$  for every input length  $n \geq 1$ :

$$\text{FUB}_s \triangleq \forall N \forall n \exists F (n = |N| \geq 1 \wedge \text{Size}(F) < s(n) \wedge \forall x (|x| \leq n \rightarrow \text{Eval}(F, x) = \oplus_n^0(x)).$$

516 **Theorem 14.** *Let  $s(n) \triangleq 4n^2$ . Then  $\text{PV}_1 \vdash \text{FUB}_s$ .*

517 *Proof.*  $\text{FUB}_s$  is a  $\forall\Sigma_2^b$  sentence and our intended theory is  $\text{PV}_1$ . In order to implement some inductive  
518 proofs, it will be helpful to reduce the complexity of the formula. For this, we introduce a new polynomial-  
519 time function,  $\text{ParForm}(1^n)$ , which generates the desired formula that computes the parity function on  $n$   
520 bits. Since it is a polynomial-time function, there is a symbol for it in  $\text{PV}$  and we can use it in the new  
521 formalization:

$$\text{FUB}'_s \triangleq \forall N \forall n (n = |N| \geq 1 \wedge \text{Size}(\text{ParForm}(1^n)) < s(n) \wedge \forall x (|x| \leq n \rightarrow \text{Eval}(\text{ParForm}(1^n), x) = \oplus_n^0(x)).$$

522 It is immediate that  $\text{FUB}'_s \Rightarrow \text{FUB}_s$ , thus we focus on proving  $\text{FUB}'_s$ . We continue with the following steps:

523 1. We prove an upper bound of  $n^2$  for the formulas calculating the parity function and its negation, when  
 524  $n$  is a power of 2.

525 2. We use this construction to derive the  $4n^2$  upper bound for any  $n$ .

526 Next, we define a polynomial-time algorithm  $\text{Par}(1^n)$  which computes a formula that calculates the  
 527 parity function on  $n$  bits and a formula that calculates the negation of the parity function on  $n$  bits, if  $n$  is a  
 528 power of 2.

**Input:**  $1^n$  for some  $n \geq 1$ .

1 Let  $k \triangleq \lfloor n/2 \rfloor$ . If  $n \neq 2^k$  ( $n$  is not a power of 2), then **return** “error”;  
 //  $F$  will compute the parity function, while  $\overline{F}$  will compute its  
 negation

2 **if**  $k = 0$  **then**

3 | Define  $F$  to be the formula with one leaf  $x_1$  and  $\overline{F}$  to be the formula with one leaf  $\neg x_1$ .

4 **else if**  $k \geq 1$  **then**

529 | // Construct a pair  $(F, \overline{F})$  of formulas on input bits  $x_1, \dots, x_{2^k}$  as  
 follows:

5 | Let  $(F_1, \overline{F}_1) \triangleq \text{Par}(1^{2^{k-1}})$ , and define a corresponding pair  $(F_2, \overline{F}_2)$ :

6 | In  $F_2$  and  $\overline{F}_2$ , relabel the leaves by putting  $x_{2^{k-1}+i}$  instead of  $x_i$  for every  $i = 1, \dots, 2^{k-1}$ ;

7 | Now let  $F \triangleq (F_1 \vee F_2) \wedge (\overline{F}_1 \vee \overline{F}_2)$  and  $\overline{F} \triangleq (F_1 \wedge F_2) \vee (\overline{F}_1 \wedge \overline{F}_2)$ .

8 **end**

9 **return**  $(F, \overline{F})$ .

**Algorithm 4:**  $\text{Par}(1^n)$  outputs Boolean formulas for  $\oplus_n^0$  and  $\oplus_n^1$  when  $n$  is a power of 2.

530 **Lemma 15.** *If  $n$  is a power of 2, the algorithm  $\text{Par}(1^n)$  correctly outputs two formulas  $(F, \overline{F})$  of size  $n^2$   
 531 which calculate the parity function and its negation, provably in  $\text{S}_2^1(\text{PV})$ .*

532 *Proof.* We split the proof of the correctness for the algorithm  $\text{Par}(1^n)$  into 3 properties:

- 533 1.  $\phi_1(n) \triangleq F, \overline{F} \in \text{VALIDFORM}(n)$ , where  $\text{VALIDFORM}(n)$  is the set of formulas on  $n$  variables;
- 534 2.  $\phi_2(n) \triangleq \text{Size}(F) = \text{Size}(\overline{F}) = n^2$ ;
- 535 3.  $\phi_3(n) \triangleq \forall x \ |x| \leq n \rightarrow \text{Eval}(F, x) = \oplus_n^0(x) \wedge \text{Eval}(\overline{F}, x) = \oplus_n^1(x)$ .

536 For now we only care about the case that  $n$  is a power of 2, so we prove these properties conditionally  
 537 (equivalently we prove  $(n = (n-1)\#1) \rightarrow \phi(n)$ ).<sup>7</sup> That is why it suffices to use polynomial induction on  
 538  $n$ , which is available in  $\text{S}_2^1$ , since our formulas are at most  $\Pi_1^b$ .

539 We skip the proof of  $\phi_1$ , which is proven by simple induction as below, using the fact that if  $F_1, F_2$  are  
 540 formulas then  $F_1 \wedge F_2$  and  $F_1 \vee F_2$  are also formulas.

**Property 2:**  $\text{S}_2^1 \vdash \phi_2(n)$ . For the base case,  $\phi_2(1)$ , we have  $k = 0$ , which means that the output  $(F, \overline{F}) \triangleq \text{Par}(1^1)$  will be two formulas with one leaf each, hence

$$\text{Size}(F) = \text{Size}(\overline{F}) = 1.$$

<sup>7</sup>It is easy to check that this is true if and only if  $n$  is a power of 2.

For the induction step, we need  $S_2^1 \vdash \forall n \phi_2(\lfloor n/2 \rfloor) \rightarrow \phi_2(n)$ . If  $n$  is not a power of 2, then the statement is true by default. In the case of  $n$  being a power of 2, we fix  $k = \lfloor n - 1 \rfloor$  and we want to prove equivalently:

$$S_2^1 \vdash \phi_2(2^{k-1}) \rightarrow \phi_2(2^k).$$

541 Assume that  $\phi_2(2^{k-1}) \equiv \phi_2(n/2)$  holds. From Line 8 we have that

$$F = (F_1 \vee F_2) \wedge (\overline{F_1} \vee \overline{F_2}) \text{ and } \overline{F} = (F_1 \wedge F_2) \vee (\overline{F_1} \wedge \overline{F_2}), \quad (10)$$

where  $(F_1, \overline{F_1})$  and  $(F_2, \overline{F_2})$  are copies of  $\text{Par}(1^{n/2})$ . From the induction hypothesis, this means that  $\text{Size}(F_1) = \text{Size}(\overline{F_1}) = \text{Size}(F_2) = \text{Size}(\overline{F_2}) = (n/2)^2 = 2^{2(k-1)}$ . Therefore, from (Equation (10)) and the properties of the function  $\text{Size}$ , we get

$$\text{Size}(F) = \text{Size}(F_1) + \text{Size}(\overline{F_1}) + \text{Size}(F_2) + \text{Size}(\overline{F_2}) = 4 \cdot 2^{2(k-1)} = 2^{2k} = n^2.$$

542 Similarly for  $\overline{F}$ , which means that  $\phi_2(2^k) \equiv \phi_2(n)$  holds. This completes the proof of the induction for  
543  $\phi_2$ .

544 **Property 3:**  $S_2^1 \vdash \phi_3(n)$ . Here the base case is trivial: for  $F \triangleq x_1$  and  $x \in \{0, 1\}$ , then  $\text{Eval}(F, x) = x =$   
545  $\oplus_1^0(x)$ . Similarly for  $\overline{F}$ .

For the induction step, we assume as above that  $n = 2^k$  and we want to prove:

$$S_2^1 \vdash \phi_3(2^{k-1}) \rightarrow \phi_3(2^k).$$

546 We assume that  $\phi_2(2^{k-1}) \equiv \phi_2(n/2)$  holds and we write  $F$  in the form

$$F = (F_1 \vee F_2) \wedge (\overline{F_1} \vee \overline{F_2}) \text{ and } \overline{F} = (F_1 \wedge F_2) \vee (\overline{F_1} \wedge \overline{F_2}),$$

where  $(F_1, \overline{F_1})$  and  $(F_2, \overline{F_2})$  are copies of  $\text{Par}(1^{n/2})$ . Therefore, instead of  $\text{Eval}(F, x)$ , we can calculate

$$\text{Eval}((F_1 \vee F_2) \wedge (\overline{F_1} \vee \overline{F_2}), x).$$

547 We need to prove that  $\text{Eval}(F, x) = \oplus_n^0(x)$  for all  $x$  with  $|x| \leq n$ . So, taking one such  $x$  we can split  
548 its binary representation into two parts  $x_1, x_2$  with lengths  $|x_1|, |x_2| \leq n/2$ , such that  $x = (x_2 x_1)_b =$   
549  $x_1 + 2^{n/2} x_2$ .

550 The input to subformulas  $F_2, \overline{F_2}$  from the definition are the bits  $x_{2^{k-1}+i}$  for  $i = 1, \dots, 2^{k-1}$ , which  
551 means that their input is  $x_2$ . Similarly, the input to subformulas  $F_1, \overline{F_1}$  is  $x_1$ . Hence, we can define

$$\begin{aligned} b_1 &\triangleq \text{Eval}(F_1, x_1) & b_3 &\triangleq \text{Eval}(\overline{F_1}, x_1) \\ b_2 &\triangleq \text{Eval}(F_2, x_2) & b_4 &\triangleq \text{Eval}(\overline{F_2}, x_2) \end{aligned}$$

552 From the properties of the evaluation function and the form of  $F$ , we can prove in  $S_2^1$  that  $\text{Eval}(F, x) =$   
553  $(b_1 \vee b_2) \wedge (b_3 \vee b_4)$ , where the symbols  $\vee, \wedge$  are used as Boolean symbols here.

554 However, since  $|x_1|, |x_2| \leq n/2$  and  $(F_1, \overline{F_1}) = (F_2, \overline{F_2}) = \text{Par}(1^{n/2})$ , from the induction hypothesis  
555 we get that

$$\begin{aligned} b_1 &= \oplus^0(x_1) & b_3 &= \oplus^1(x_1) = 1 - b_1 \\ b_2 &= \oplus^0(x_2) & b_4 &= \oplus^1(x_2) = 1 - b_2 \end{aligned}$$

Next, it is easy to prove by checking all the 4 cases that

$$\forall b_1, b_2 \in \{0, 1\} (b_1 \vee b_2) \wedge ((1 - b_1) \vee (1 - b_2)) = b_1 \oplus b_2,$$

and as a result, we get

$$\text{Eval}(F, x) = (\oplus^0(x_1)) \oplus (\oplus^0(x_2)) = \oplus^0(x_2x_1) = \oplus^0(x)$$

556 by the properties of the parity function. Similarly, we can prove that  $\text{Eval}(\overline{F}, x) = \oplus_n^1(x)$ , which concludes  
557 the induction.  $\square$

For the general case, we use a simple padding argument. For a number  $n$ , we can define the number

$$\tilde{n} \triangleq (n - 1)\#1.$$

This number is the least power of 2 that is greater or equal to  $n$ . It is easy to see that

$$\text{PV}_1 \vdash n \leq \tilde{n} < 2n.$$

558 If we replace  $\text{ParForm}(1^n)$  by  $\text{Par}_1(1^{\tilde{n}})$  (the first coordinate of  $\text{Par}(1^{\tilde{n}})$ ), we have by the above lemma  
559 that

- 560 1.  $\text{Size}(\text{ParForm}(1^n)) = \text{Size}(\text{Par}_1(1^{\tilde{n}})) = \tilde{n}^2 < (2n)^2 = s(n)$ .
- 561 2. For all  $x$  with  $|x| \leq n$ , we have  $|x| \leq \tilde{n}$ , which by the lemma gives us  $\text{Eval}(\text{ParForm}(1^n), x) =$   
562  $\text{Eval}(\text{Par}_1(1^{\tilde{n}}), x) = \oplus_n^0(x)$ . Since  $|x| \leq n$ , we also have  $\oplus_n^0(x) = \oplus_n^0(x)$ . Consequently, we have  
563  $\text{Eval}(\text{ParForm}(1^n), x) = \oplus_n^0(x)$ .

564 These two together show that  $\text{PV}_1 \vdash \text{FUB}'_s$  and the proof is complete.  $\square$

### 565 4.3 Formula Size Hierarchy

566 In this section, we provide the proof of Theorem 3.

**Theorem 16** (Theorem 3). *Consider rationals  $a > 2$  and  $b = 3/2$ , and let  $n_0$  be a large enough positive integer. Then*

$$\text{PV}_1 \vdash \text{FSH}[a, b, n_0].$$

567 *Proof.* We combine the results of Section 4.1 and Section 4.2. We argue in  $\text{PV}_1$ . From Theorem 11, we get  
568 that

$$\forall n \in \text{Log} \forall F \in \text{FORMULA}[n^{3/2}] \exists x (|x| \leq n \wedge F(x) \neq \oplus_n(x)), \quad (11)$$

569 and from Theorem 14, we have that

$$\forall n \in \text{Log} \exists G \in \text{FORMULA}[4n^2] \forall x (|x| \leq n \rightarrow F(x) = \oplus_n(x)).$$

570 We can eliminate the constant 4 from the latter using that  $a > 2$  and choosing a large enough  $n_0$ , such that  
571 for every  $n \geq n_0$ ,  $n^a \geq 4n^2$  (provably in  $\text{PV}_1$ ). Consequently,

$$\forall n \geq n_0 \in \text{Log} \exists G \in \text{FORMULA}[n^a] \forall x (|x| \leq n \rightarrow F(x) = \oplus_n(x)). \quad (12)$$

Finally, combining Equation (11) and Equation (12), we get that

$$\forall n \geq n_0 \in \text{Log} \exists G \in \text{FORMULA}[n^a] \forall F \in \text{FORMULA}[n^{3/2}] \exists x (|x| \leq n \wedge F(x) \neq G(x)),$$

572 which is exactly the formula size hierarchy,  $\text{FSH}[a, b, n_0]$ , for our choice of parameters  $a > 2$  and  $b =$   
573  $3/2$ .  $\square$

## References

- 574
- 575 [AB09] Sanjeev Arora and Boaz Barak. *Computational Complexity - A Modern Approach*. Cambridge  
576 University Press, 2009. 6
- 577 [Bus86] Samuel R. Buss. *Bounded Arithmetic*. Bibliopolis, 1986. 3, 7
- 578 [Bus97] Samuel R. Buss. Bounded arithmetic and propositional proof complexity. In *Logic of Com-*  
579 *putation*, pages 67–121. Springer Berlin Heidelberg, 1997. 3
- 580 [CLO24] Lijie Chen, Jiatu Li, and Igor C. Oliveira. Reverse mathematics of complexity lower bounds.  
581 In *Symposium on Foundations of Computer Science (FOCS)*, 2024. 5
- 582 [CMMW19] Lijie Chen, Dylan M. McKay, Cody D. Murray, and R. Ryan Williams. Relations and equiva-  
583 lences between circuit lower bounds and Karp-Lipton theorems. In *Computational Complexity*  
584 *Conference (CCC)*, pages 30:1–30:21, 2019. 3
- 585 [CN10] Stephen A. Cook and Phuong Nguyen. *Logical Foundations of Proof Complexity*. Cambridge  
586 University Press, 2010. 7
- 587 [Cob65] Alan Cobham. The intrinsic computational difficulty of functions. *Proc. Logic, Methodology*  
588 *and Philosophy of Science*, pages 24–30, 1965. 7
- 589 [Coo75] Stephen A. Cook. Feasibly constructive proofs and the propositional calculus (preliminary  
590 version). In *Symposium on Theory of Computing (STOC)*, pages 83–97, 1975. 3, 7
- 591 [IW97] Russell Impagliazzo and Avi Wigderson. P=BPP if E requires exponential circuits: De-  
592 randomizing the XOR lemma. In *Symposium on the Theory of Computing (STOC)*, pages  
593 220–229, 1997. 3
- 594 [Jeř04] Emil Jeřábek. Dual weak pigeonhole principle, boolean complexity, and derandomization.  
595 *Annals of Pure and Applied Logic*, 129(1-3):1–37, 2004. 3, 7
- 596 [Jeř05] Emil Jeřábek. *Weak pigeonhole principle and randomized computation*. PhD thesis, Charles  
597 University in Prague, 2005. 3, 7
- 598 [Jeř06] Emil Jeřábek. The strength of sharply bounded induction. *Mathematical Logic Quarterly*,  
599 52(6):613–624, 2006. 3, 7
- 600 [Jeř07] Emil Jeřábek. Approximate counting in bounded arithmetic. *Journal of Symbolic Logic*,  
601 72(3):959–993, 2007. 3, 7
- 602 [Juk12] Stasys Jukna. *Boolean Function Complexity: Advances and Frontiers*. Springer, 2012. 14,  
603 15, 16, 17
- 604 [KO17] Jan Krajíček and Igor C. Oliveira. Unprovability of circuit upper bounds in Cook’s theory PV.  
605 *Logical Methods in Computer Science*, 13(1), 2017. 5
- 606 [KPT91] Jan Krajíček, Pavel Pudlák, and Gaisi Takeuti. Bounded arithmetic and the polynomial hier-  
607 archy. *Annals of Pure and Applied Logic*, 52(1-2):143–153, 1991. 8



- 608 [Kra92] Jan Krajíček. No counter-example interpretation and interactive computation. In *Logic from*  
609 *Computer Science: Proceedings of a Workshop held November 13–17, 1989*, pages 287–293.  
610 Springer, 1992. 14
- 611 [Kra95] Jan Krajíček. *Bounded Arithmetic, Propositional Logic, and Complexity Theory*. Encyclope-  
612 dia of Mathematics and its Applications. Cambridge University Press, 1995. 7, 8, 16
- 613 [Kra19] Jan Krajíček. *Proof Complexity*. Encyclopedia of Mathematics and its Applications. Cam-  
614 bridge University Press, 2019. 7, 13
- 615 [Kra21] Jan Krajíček. Small circuits and dual weak PHP in the universal theory of p-time algorithms.  
616 *ACM Transactions on Computational Logic (TOCL)*, 22(2):1–4, 2021. 5
- 617 [MP20] Moritz Müller and Ján Pich. Feasibly constructive proofs of succinct weak circuit lower  
618 bounds. *Annals of Pure and Applied Logic*, 171(2), 2020. 5
- 619 [MPW02] Alexis Maciel, Toniann Pitassi, and Alan R. Woods. A new proof of the weak pigeonhole  
620 principle. *Journal of Computer and System Sciences*, 64(4):843–872, 2002. 7, 13
- 621 [Oli24] Igor C. Oliveira. Meta-mathematics of computational complexity theory. *Preprint*, 2024. 3,  
622 4, 5
- 623 [PWW88] Jeff B. Paris, A. J. Wilkie, and Alan R. Woods. Provability of the pigeonhole principle and the  
624 existence of infinitely many primes. *J. Symb. Log.*, 53(4):1235–1244, 1988. 13
- 625 [Sha49] Claude E. Shannon. The synthesis of two-terminal switching circuits. *The Bell System Tech-*  
626 *nical Journal*, 28(1):59–98, 1949. 3
- 627 [Sub61] Bella A. Subbotovskaya. Realization of linear functions by formulas using  $+$ ,  $\cdot$ ,  $-$ . In *Soviet*  
628 *Math. Dokl*, 1961. 6, 14

## 629 A Proof of the KPT Theorem for $\forall\exists\forall\exists$ Sentences

630 In order to make our results more accessible and the presentation self-contained, in this section we  
631 describe a standard model-theoretic proof of the KPT Witnessing Theorem. We restate the result below for  
632 convenience of the reader.

**Theorem 17.** *Let  $\mathbb{T}$  be a universal theory with vocabulary  $\mathcal{L}$ . Let  $\varphi$  be an open  $\mathcal{L}$ -formula, and suppose that*

$$\mathbb{T} \vdash \forall x \exists y \forall z \exists w \varphi(x, y, z, w).$$

*Then there is a finite sequence  $s_1, \dots, s_k$  of  $\mathcal{L}$ -terms such that*

$$\mathbb{T} \vdash \forall x, z_1, \dots, z_k (\psi(x, s_1(x), z_1) \vee \psi(x, s_2(x, z_1), z_2) \vee \dots \vee \psi(x, s_k(z_1, \dots, z_{k-1}), z_k)),$$

*where*

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

633 *Proof.* Let  $b, c_1, c_2, \dots$  be a list of new constants, and let  $u_1, u_2, \dots$  be an enumeration of all terms built  
 634 from the functions and constants in  $\mathcal{L}$  together with  $b, c_1, c_2, \dots$ , where the only new constants in  $u_k$  are  
 635 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  $\mathbb{T}$  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  $\mathbb{T}$  satisfies the negation of this sentence, and by the completeness theorem,

$$\mathbb{T} \vdash \psi(b, u_1, c_1) \vee \psi(b, u_2, c_2) \vee \dots \vee \psi(b, u_k, c_k) .$$

636 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  
 637 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

$$\mathbb{T} \cup \{\neg\psi(b, u_1, c_1), \neg\psi(b, u_2, c_2), \neg\psi(b, u_3, c_3), \dots\}$$

admits a model  $\mathcal{M}$ . Consequently, using the definition of  $\psi$ ,

$$\mathcal{M} \models \mathbb{T} \cup \{\forall w \neg\varphi(b, u_1, c_1, w), \forall w \neg\varphi(b, u_2, c_2, w), \dots\}$$

Let  $\mathbb{T}^+ \triangleq \mathbb{T} \cup \{\forall w \neg\varphi(b, u_1, c_1, w), \forall w \neg\varphi(b, u_2, c_2, w), \dots\}$ . Since  $\mathbb{T}$  is a universal theory and  $\varphi$  is an open formula, it follows that  $\mathbb{T}^+$  is also a universal theory. For this reason, the substructure  $\mathcal{M}'$  of  $\mathcal{M}$  consisting of the denotations of the terms  $u_1, u_2, \dots$  is also a model of  $\mathbb{T}^+$ . Now it is not hard to prove that

$$\mathcal{M}' \models \mathbb{T} + \exists x \forall y \exists z \forall w \neg\varphi(x, y, z, w) ,$$

638 which contradicts the hypothesis of the theorem and completes the proof. To see this, it is enough to show  
 639 that  $\mathcal{M}' \models \forall y \exists z \forall w \neg\varphi(b^{\mathcal{M}'}, y, z, w)$ . Given an arbitrary element  $m$  in  $\mathcal{M}'$ , by construction of  $\mathcal{M}'$ , there  
 640 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  $\mathbb{T}^+$ , which includes  
 641 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  
 642  $\mathcal{M}' \models \forall y \exists z \forall w \neg\varphi(b^{\mathcal{M}'}, y, z, w)$ .  $\square$