



# Finding Bugs in Short Proofs: The Metamathematics of Resolution Lower Bounds

Jiawei Li  
UT Austin

[davidlee@cs.utexas.edu](mailto:davidlee@cs.utexas.edu)

Yuhao Li  
Columbia University

[yuhaoli@cs.columbia.edu](mailto:yuhaoli@cs.columbia.edu)

Hanlin Ren  
University of Oxford

[hanlin.ren@cs.ox.ac.uk](mailto:hanlin.ren@cs.ox.ac.uk)

## Abstract

We study the *refuter* problems for proof complexity lower bounds. Suppose  $\varphi$  is a hard tautology that does not admit any length- $s$  proof in some proof system  $P$ . In the corresponding refuter problem, we are given (query access to) a purported length- $s$  proof  $\pi$  in  $P$  that claims to have proved  $\varphi$ , and our goal is to find an invalid derivation step within  $\pi$ . As suggested by witnessing theorems in bounded arithmetic, the *computational complexity* of these refuter problems is closely tied to the *metamathematics* of the underlying lower bounds.

We focus on refuter problems corresponding to lower bounds for *resolution*, which is arguably the single most studied system in proof complexity. As a warm-up, we show that many refuter problems for resolution *width* lower bounds are PLS-complete. To capture the complexity of refuter problems for resolution *size* lower bounds, we introduce a new class  $\text{rwPHP(PLS)}$  in decision-tree TFNP, which can be seen as a randomized version of PLS.

- We show that the refuter problems for many resolution size lower bounds can be solved in  $\text{rwPHP(PLS)}$ , including the classic lower bound of Haken [TCS, 1985] for the pigeonhole principle. More generally, we identify a common proof technique that we call “random restriction + width lower bound”, and present strong evidence that resolution lower bounds proved by this technique typically have refuter problems in  $\text{rwPHP(PLS)}$ .
- We then show that the refuter problem for *any* resolution size lower bound is  $\text{rwPHP(PLS)}$ -hard, thereby demonstrating that the  $\text{rwPHP(PLS)}$  upper bound mentioned above is tight. Informally speaking, this means that “ $\text{rwPHP(PLS)}$ -reasoning” is *necessary* for proving *all* resolution size lower bounds.

Interpreted in bounded arithmetic, our results show that the theory  $\text{T}_2^1(\alpha) + \text{dwPHP(PV}(\alpha))$  characterizes the “reasoning power” required to prove (the “easiest”) resolution size lower bounds.

As a corollary, we obtain surprisingly efficient proofs of resolution lower bounds. In particular, we show that many resolution size lower bounds can be proved in low-width *random resolution* [Pudlák–Thapen, CCC’17].

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	More Background . . . . .	3
1.2	Our Settings . . . . .	4
1.2.1	Refuter Problems for Resolution Lower Bounds . . . . .	5
1.2.2	Retraction Weak Pigeonhole Principles . . . . .	6
1.3	Our Results . . . . .	7
1.3.1	Refuters for Resolution Width Lower Bounds . . . . .	8
1.3.2	Refuters for Resolution Size Lower Bounds . . . . .	9
1.3.3	Applications . . . . .	10
1.4	Discussions, Speculations, and Future Directions . . . . .	12
1.5	Further Related Works . . . . .	13
<b>2</b>	<b>Technical Overview</b>	<b>14</b>
2.1	Refuter Problems in rwPHP(PLS) . . . . .	14
2.2	Refuter Problems are rwPHP(PLS)-Hard . . . . .	16
<b>3</b>	<b>Preliminaries</b>	<b>17</b>
3.1	Pigeonhole Principle . . . . .	18
3.2	Decision Tree TFNP . . . . .	18
3.2.1	Connection to Proof Complexity . . . . .	19
3.3	Bounded Arithmetic . . . . .	20
3.4	Refuter Problems for Resolution Lower Bounds . . . . .	22
3.5	$\mathcal{P}$ -Retraction Weak Pigeonhole Principle . . . . .	23
3.5.1	Witnessing for $T_2^1 + \text{dwPHP(PV)}$ . . . . .	24
<b>4</b>	<b>Refuters for the Pigeonhole Principle</b>	<b>25</b>
4.1	Refuters for Narrow Resolution Proofs . . . . .	25
4.2	Refuters for Short Resolution Proofs . . . . .	29
<b>5</b>	<b>Hardness of Refuting Resolution Proofs</b>	<b>33</b>
5.1	Hardness of Refuting Narrow Resolution Proofs . . . . .	33
5.2	Hardness of Refuting Short Resolution Proofs . . . . .	34
<b>6</b>	<b>Refuters for Other Formulas</b>	<b>39</b>
6.1	Universal Refuters for <i>Every</i> Narrow Resolution Proof . . . . .	39
6.2	Refuters for XOR-Lifted Lower Bounds . . . . .	40
6.3	Refuters for Tseitin Formulas . . . . .	42
6.4	Refuters for Random $k$ -CNFs . . . . .	46
6.5	Open Problems: What We <i>Failed</i> to Formalize . . . . .	50
<b>7</b>	<b>Applications</b>	<b>50</b>
7.1	Proof Complexity of Proof Complexity Lower Bounds . . . . .	50
7.2	Complexity of Black-Box TFNP Separations . . . . .	53
7.2.1	Black-Box TFNP Refuters and its Properties . . . . .	53
7.2.2	Refuter for Separating from PLS . . . . .	55
	<b>References</b>	<b>57</b>
<b>A</b>	<b>Amplification for rwPHP(<math>\mathcal{P}</math>)</b>	<b>63</b>
<b>B</b>	<b>Comparing REFUTER with WRONGPROOF(Res)</b>	<b>66</b>
<b>C</b>	<b>Prover-Delayer Games, PLS, and the Proof of Lemma 7.10</b>	<b>68</b>
C.1	From PLS to Resolution using Prover-Delayer Game . . . . .	68
C.2	Proof of Lemma 7.10 . . . . .	70

# 1 Introduction

One of the earliest lower bounds in proof complexity was Haken’s landmark result [Hak85] that the pigeonhole principle requires exponential-size proofs in the resolution proof system. Since then, proof complexity has become a vibrant research area with substantial progress in establishing lower bounds for various proof systems, as well as the development of a wide range of lower bound techniques. However, despite decades of efforts, proving nontrivial lower bounds for stronger systems, such as Frege and Extended Frege, remains elusive. It is widely believed that proving lower bounds for Extended Frege is “beyond our current techniques”<sup>1</sup>, but what does this even mean? How much, and in which directions, must our techniques expand, to enable us to prove lower bounds for stronger proof systems? These questions call for a study of the *metamathematical* difficulty of proving lower bounds in proof complexity (see, e.g., [PS19, ST21]).

Inspired by recent works on the reverse mathematics of *circuit lower bounds* [CJSW24, Kor22, CTW23, CLO24], we propose investigating the metamathematics of proof complexity lower bound through the computational lens of their *refuter* problem. To illustrate, consider the following total search problem: suppose we are given a resolution proof  $\Pi$  that claims to prove the pigeonhole principle, yet its length is shorter than the lower bound established in [Hak85]. By Haken’s result,  $\Pi$  cannot be a valid resolution proof; it must contain an invalid derivation. The goal of the search problem is to locate such an error. We refer to this total search problem as the “refuter problem”<sup>2</sup> corresponding to Haken’s lower bound:

**Problem 1.1** (Informal). Given (query access to) a subexponential-size resolution proof  $\Pi$  that claims to be a proof of the pigeonhole principle, find an invalid derivation in  $\Pi$ .

For any proof complexity lower bound of the form “the tautology  $\phi$  requires proof length greater than  $s$  in the proof system  $P$ ”, we can define an associated search problem: Given a purported  $P$ -proof  $\Pi$  of  $\phi$  with length at most  $s$ , find an invalid derivation in  $\Pi$ . With the appropriate formalization (see Section 1.2), these refuter problems are NP *search problems* and are *total* if and only if their underlying lower bounds hold. Therefore, their computational complexity can be studied using the theory of TFNP [MP91]. As elaborated in Section 1.1, the complexity of refuter problems reflects the metamathematical difficulty of proving the corresponding lower bounds, thereby providing a purely computational framework for analyzing the difficulty of proving such lower bounds.

In this paper, we initiate the study of refuter problems in proof complexity and take the first step by studying these problems for *resolution* lower bounds. Resolution serves as a natural first step for exploring the metamathematics of proof lower bounds for two main reasons:

- (i) First, resolution is a well-studied proof system, largely due to its fundamental connections to SAT-solving and automated theorem provers [DP60, DLL62]. Krajíček even estimates that “there are perhaps more papers published about proof complexity of resolution than about all remaining proof complexity topics combined” [Kra19, Chapter 13].
- (ii) Second, significant progress has already been made in proving lower bounds against resolution [Hak85, Urq87, CS88, BP96, BW01], suggesting that investigating the metamathematics of resolution lower bounds is a promising avenue.

We study several important resolution lower bounds, including those for the pigeonhole principle [Hak85, BP96], Tseitin tautologies [Urq87, Sch97], and random CNF formulas [CS88]. In our study, we introduce

---

<sup>1</sup>This belief is partly supported by the intuition that proving strong circuit lower bounds (e.g.,  $\text{NP} \not\subseteq \text{P}/\text{poly}$ ) seems to be a prerequisite for proving strong proof complexity lower bounds (e.g., for Extended Frege) [Raz15]. However, formalizing such connections has proven challenging [PS26, AKPS24].

<sup>2</sup>This term is adopted from [CTW23], as will be discussed later.

a new syntactic subclass of decision-tree TFNP, denoted as rwPHP(PLS), which can be thought of as a randomized version of PLS, residing slightly above PLS in the TFNP hierarchy<sup>3</sup>.

At a high level, we show that resolution *width* lower bounds are captured by PLS, while resolution *size* lower bounds are captured by rwPHP(PLS).

**Theorem 1.2** (Main Result; Informal). *The refuter problems corresponding to the*

1. *resolution width lower bounds for the pigeonhole principle and Tseitin tautologies are PLS-complete;*
2. *resolution size lower bounds for the pigeonhole principle (Problem 1.1), Tseitin tautologies, and random CNF formulas are rwPHP(PLS)-complete.*

Our results are more comprehensive than those stated in Theorem 1.2, and we defer the full formal presentation to Section 1.3. Here, we highlight a few key insights:

- **Characterizing a common proof technique:** All the aforementioned resolution size lower bound proofs share a common strategy, which we call “random restrictions + width lower bounds”. It is often the case that resolution lower bounds proven using this strategy have refuter problems in rwPHP(PLS)<sup>4</sup>. This implies that “rwPHP(PLS)-reasoning” is sufficient for implementing one of the most commonly employed techniques for proving resolution lower bounds.
- **Minimum reasoning for resolution lower bounds:** Complementing the above, we prove that for *any* family of hard tautologies for resolution, the corresponding refuter problem (for size lower bound) is rwPHP(PLS)-hard. This establishes that our rwPHP(PLS) upper bound is indeed tight. Notably, the hardness proof does not rely on the hard tautology being the pigeonhole principle. Consequently, this result carries an intriguing metamathematical implication: “rwPHP(PLS)-reasoning” is necessary for proving *any* resolution lower bound.
- **Consequences in bounded arithmetic:** Theorem 1.2 can also be interpreted as conservativeness results showing that a certain fragment of (relativized) bounded arithmetic,  $\mathcal{T}_{\text{Res}} := \text{T}_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha))$ , “captures” the minimum reasoning required for proving resolution size lower bounds. More precisely,  $\mathcal{T}_{\text{Res}}$  is powerful enough to formalize many resolution lower bounds proved in the literature, including Haken’s seminal lower bound for PHP [Hak85], while at the same time, it is necessary for proving *any* resolution lower bound. An interesting takeaway of this result is that it is consistent with every theory weaker than  $\mathcal{T}_{\text{Res}}$  that *resolution is p-bounded* (i.e., proves every tautology in polynomial size).

We find the *existence* of such a theory quite insightful. It is natural to speculate that there is a very powerful theory  $\mathcal{T}_{\text{EF}}$  such that strong proof systems like Extended Frege “appears p-bounded” to every theory weaker than  $\mathcal{T}_{\text{EF}}$ . If true, this speculation would provide a strong barrier result in proof complexity.

- **Surprisingly efficient proofs for resolution lower bounds:** Finally, translating our results into proof complexity, we exhibit surprisingly efficient proofs of resolution lower bounds. For instance, we show that low-width *random resolution* [PT19] can prove many resolution size lower bounds (including Haken’s lower bound).<sup>5</sup> Previously, Cook and Pitassi [CP90] proved Haken’s lower bound in the theory  $\text{IPV}^\omega$ , which can be interpreted as intuitionistic reasoning using polynomial-time concepts over the purported resolution proof. Our results suggest that  $\text{AC}^0$  concepts (and indeed much weaker ones) already suffice to prove the same lower bound.<sup>6</sup>

<sup>3</sup>The formal definition and properties of rwPHP(PLS) are presented in Section 3.5.

<sup>4</sup>An interesting exception is the general size-width tradeoff by Ben-Sasson and Wigderson [BW01]; see Section 6.5 for further discussions.

<sup>5</sup>Although random resolution is not a Cook–Reckhow proof system unless  $\text{P} = \text{NP}$  [PT19], it is possible to define a fragment of random resolution that is Cook–Reckhow and that proves the aforementioned resolution size lower bounds.

<sup>6</sup>Our formalization is different from [CP90]. Informally speaking, [CP90] uses “*polynomial-time* reasoning” to handle a *polynomial-size* proof, while we use “*PH-reasoning*” to handle an *exponentially long* proof.

## 1.1 More Background

We now provide further background on bounded reverse mathematics, refuter problems, and the theory of TFNP to justify our methodology: the *metamathematics* of the proof complexity lower bounds can — and indeed *should* — be understood through the *computational complexity* of their associated refuter problems within TFNP.

**Bounded reverse mathematics.** Reverse mathematics explores, for each mathematical theorem of interest, the minimal theory required to prove it. In bounded reverse mathematics [Coo07, Ngu08, CN10], the theories considered come from *bounded arithmetic*, which (roughly speaking) are logical theories formalizing the idea of “reasoning within a complexity class  $C$ ”. The link between these logical theories and complexity classes makes bounded arithmetic, and hence bounded reverse mathematics, an effective framework for studying the metamathematics of complexity theory.

Indeed, there has been a long history of studying the (un)provability of lower bounds in the context of bounded arithmetic: In 1989, Krajíček and Pudlák investigated the unprovability of proof lower bounds [KP89], while Razborov studied the unprovability of circuit lower bounds in 1995 [Raz95a, Raz95b]. Notably, many lower bounds for weak circuit classes and proof systems can be formalized in weak theories [Raz95a, CP90, MP20], while some strong lower bounds are unprovable within them [KP89, Raz95b, Kra97, Kra11b, Pic15, PS21, LO23, CLO25].

We take a different perspective from the aforementioned line of work: rather than asking whether lower bounds are provable in certain theories, our goal is to *characterize* the exact reasoning power required to prove these lower bounds. That is, we seek to identify the *minimal* theory  $\mathcal{T}$  that can prove the given lower bound and to establish the minimality of  $\mathcal{T}$  by showing that the axioms used in the proof are indeed *necessary*. The necessity of axioms, i.e., deriving the axiom back from the theorem, is called a *reversal* in reverse mathematics.

**Example 1.3.** Recently, Chen, Li, and Oliveira [CLO24] presented several notable reversals related to complexity lower bounds. In their work, they establish that variants of weak pigeonhole principles are *necessary and sufficient* for proving various classical lower bounds. For instance, the fact that one-tape Turing machines require  $\Omega(n^2)$  time to recognize palindromes [Maa84] can be proved using the *weak pigeonhole principle*; Moreover, [CLO24, Theorem 4.9] demonstrates a reversal, proving that this lower bound is, in fact, *equivalent* to the weak pigeonhole principle. The work in [CLO24] serves as one of the main inspirations of this paper.

**Refuter problems.** To investigate the metamathematics of a lower bound statement, we first write down the statement in forall-exists form:

- Circuit lower bounds: Let  $L$  be a hard language and  $s$  be a size lower bound for  $L$ . The lower bound statement expresses that *for every* circuit  $C$  of size  $s$ , *there exists* an input  $x$  such that  $L(x) \neq C(x)$ .
- Proof lower bounds: Let  $\phi$  be a tautology that is hard for some proof system  $P$ , and  $s$  be a size lower bound for  $\phi$ . The lower bound statement expresses that *for every* purported  $P$ -proof  $\Pi$  of size  $s$ , *there exists* an invalid derivation step in  $\Pi$ .

In general, a statement

$$\forall x \exists y V(x, y) \tag{1}$$

would define a total search problem of finding a valid  $y$  given  $x$  such that  $V(x, y)$  holds; note that the statement is *true* if and only if the search problem is *total*. In our case, is a total search problem **Problem 1.1** precisely because of the Resolution lower bounds against the pigeonhole principle proved by Haken and others [Hak85, BP96, BW01].

This correspondence can be formally justified by the *witnessing theorems* in bounded arithmetic. A witnessing theorem for a theory  $\mathcal{T}$  links it to a syntactic subclass  $C_{\mathcal{T}}$  of TFNP, and the theorem states

that if (1) is provable in  $\mathcal{T}$ , then the corresponding (total) search problem lies in the class  $C_{\mathcal{T}}$ .<sup>7</sup> For instance, Buss’s witnessing theorem [Bus85] states that if (1) is provable in  $S_2^1$ , then the corresponding total search problem can be solved in polynomial time; Buss and Krajíček [BK94] showed that if (1) is provable in  $T_2^1$ , then the corresponding total search problem is solvable in PLS (polynomial local search).

Search problems corresponding to circuit lower bounds have already been studied in the literature [GST07, Pic15, CJSW24, Kor22, CTW23, CLO24] and are termed “refuter problems” in [CTW23]. We adopt this terminology and refer to the search problems associated with proof lower bounds as “refuter problems” as well.<sup>8</sup>

**Total search problems in NP.** The above discussion suggests that the metamathematics of lower bounds can be understood through the computational complexity of their refuter problems. Since these problems are total search problems in NP (as long as the lower bounds are true), it is natural to adopt the methodology of TFNP while studying their complexity.

What is the “methodology of TFNP”? Since the seminal work of Megiddo and Papadimitriou [MP91], problems in TFNP have been categorized based on their *proof of totality*. For instance, the class PLS captures NP search problems whose totality is provable from the principle “every DAG has a sink” [JPY88], while the class PPAD captures problems whose totality is provable from “every DAG with an unbalanced node has another one” [Pap94]. Moreover, *completeness* results play the same role as reversals in bounded reverse mathematics. For example, a pivotal result in this direction is the PPAD-completeness of finding a Nash equilibrium in two-player games [CDT09, DGP09]. This result carries an intriguing metamathematical interpretation: Topological arguments (specifically, Brouwer’s fixed point theorem [Bro11]) or methods akin to it are *unavoidable* for proving the existence of Nash equilibrium [Nas51], which stands in stark contrast to the linear programming duality methods used for zero-sum games [VS23].

The attentive reader may have already noticed that the above methodology shares a close resemblance to (bounded) reverse mathematics. This similarity can indeed be formally justified by the witnessing theorems mentioned earlier. (Another formal justification is that provability in (universal variants of) bounded arithmetic is equivalent to reducibility in TFNP; see, e.g., [Mül21, Proposition 3.4].) While reading this paper, it is useful to remember that all TFNP results established here can be translated into results in bounded arithmetic and vice versa, conveying the same underlying conceptual message.

## 1.2 Our Settings

Before explaining our results, we first discuss the setting of (decision tree) TFNP and (relativized) bounded arithmetic in which our results take place. This sub-section is of preliminary nature, but we recommend reading (i.e., not skipping) it before proceeding to our results in Section 1.3. In particular, this sub-section introduces our formalization of lower bounds, which is different from previous works [CP90, MP20] as the purported resolution proof is represented as an *exponentially-long second-order* object.

We consider TFNP problems in the *decision tree* model (TFNP<sup>dt</sup>); this model is sometimes called “type-2 TFNP problems” [BCE<sup>+</sup>98] when the decision trees are uniform. In this model, we are given an input  $x$  of length  $N$  and we think of *decision trees of polylog( $N$ ) depth* as “efficient”. Each possible solution  $o$  can be represented by polylog( $N$ ) bits, and there is an efficient procedure  $\phi(x, o)$  that verifies whether  $o$  is a valid solution for  $x$ . (That is, given the purported solution  $o$ ,  $\phi(x, o)$  makes only polylog( $N$ ) queries to  $x$ .) The goal is, of course, to find a solution  $o$  such that  $\phi(x, o)$  holds.

TFNP<sup>dt</sup> corresponds to *relativized* bounded arithmetic where a new predicate  $\alpha$  is added into the language. The predicate  $\alpha$  is intuitively treated as an oracle (or an exponentially-long input). For example, PV( $\alpha$ ) captures reasoning using  $P^\alpha$ -concepts, i.e., *uniform and efficient decision trees over  $\alpha$* .

<sup>7</sup>This requires (1) to be a “ $\forall\Sigma_1^b$ -sentence”, meaning that  $|x|$  and  $|y|$  are polynomially related and  $V(x, y)$  is a deterministic polynomial-time relation.

<sup>8</sup>In fact, [CTW23] called these problems “*refutation* problems”. We choose to use “*refuter* problems” to avoid confusion with the term “*refutation*” in proof complexity, which usually refers to a proof showing that a formula is unsatisfiable.

*Remark 1* (Type-1 vs. Type-2 TFNP Problems). In the literature, it is common to define a type-1 TFNP problem in terms of *succinct encodings* of exponentially large objects. For example, a possible definition of a PLS-complete problem is as follows: Given a “neighborhood” circuit  $C : \{0, 1\}^n \rightarrow \{0, 1\}^{\text{poly}(n)}$  and a “potential function” circuit  $V : \{0, 1\}^n \rightarrow \{0, 1\}^{\text{poly}(n)}$  that together encode a DAG on  $2^n$  nodes, and also an active node (i.e., a node with non-zero out-degree), find a sink of this graph (i.e., a node with non-zero in-degree and zero out-degree). In contrast, the TFNP<sup>dt</sup> / type-2 TFNP problems that we consider simply treat  $C$  and  $V$  as oracles.

Any separation of type-2 TFNP problems implies a separation of type-1 TFNP problems *in a relativized world* [BCE+98]. For example,  $\text{PLS}^{\text{dt}} \not\subseteq \text{PPA}^{\text{dt}}$  implies an oracle  $O$  under which  $\text{PLS}^O \not\subseteq \text{PPA}^O$ .

### 1.2.1 Refuter Problems for Resolution Lower Bounds

This subsection formalizes the refuter problem for resolution lower bounds as a TFNP<sup>dt</sup> problem. We assume familiarity with the resolution proof system. In resolution, every line is a *clause* (i.e., the disjunction of literals) and the only inference rule is the *resolution rule*:

$$\frac{C \vee \ell \quad D \vee \bar{\ell}}{C \vee D},$$

where  $C, D$  are clauses and  $\ell$  is a literal. Sometimes, we will also allow the *weakening* rule that replaces a clause with a consequence of it:

$$\frac{C}{C \vee D}.$$

The *size* of a resolution proof is the number of lines (i.e., clauses) in it. The *width* of a resolution proof is the maximum width of any clause in it, where the *width* of a clause is the number of literals in the clause. Basics about resolution can be found in any textbook on proof complexity, e.g., [Kra19, Section 5].

**Size lower bounds for resolution.** Let  $F$  be a tautology<sup>9</sup> that is *exponentially-hard* for resolution. For example, take  $F$  to be the pigeonhole principle which does not have  $c^n$ -size resolution proofs for some absolute constant  $c > 1$  [Hak85]. The refuter problem, which we denote as

$$\text{REFUTER}(s(F \vdash_{\text{Res}} \perp) \leq c^n),$$

is defined as follows. The input  $\Pi$  is a purported length- $c^n$  resolution proof of  $F$  represented as a list of  $c^n$  *nodes*, where each node consists of a clause in the resolution proof and the predecessors of this clause. (For example, if the clause in node  $i$  is resolved from the clauses in node  $j$  and node  $k$ , then the predecessor information would contain two integers  $(j, k)$ .) A *valid solution* would be the index of any node  $i \in [c^n]$  whose derivation is illegal: denoting  $C_i$  the clause in node  $i$ , then there do not exist clauses  $C, D$  and a literal  $\ell$  such that

$$C_i = C \vee D, C_j = C \vee \ell, C_k = D \vee \bar{\ell}.$$

A more formal definition can be found in Section 3.4.

By Haken’s lower bound mentioned above [Hak85], every purported resolution proof of length  $c^n$  must contain an illegal derivation, thus the above problem is *total*. Let  $N := c^n \text{poly}(n)$  denote the bit-length of the input resolution proof, then each node can be described in  $\text{poly}(n) \leq \text{polylog}(N)$  bits, hence there is an efficient decision tree that verifies whether a node  $i$  is illegal and the above refuter problem is indeed in TFNP<sup>dt</sup>.

We can also formalize resolution lower bounds in relativized bounded arithmetic as follows. We add a new symbol  $\alpha$  into our language that encodes a length- $c^n$  resolution proof, i.e., for each  $i \in [c^n]$ ,  $\alpha(i)$

<sup>9</sup>A DNF  $D$  is a *tautology* if and only if the corresponding CNF  $\neg D$  is a *contradiction*. A *proof* of  $D$  being a tautology is a *refutation* of  $\neg D$  being a contradiction. For convenience, we will use the terms “tautology/proof” and “contradiction/refutation” interchangeably.

is the  $i$ -th bit of the proof. Fixing a hard tautology  $F$ , let  $\text{mistake}_F(n, \alpha, i)$  be a  $\text{PV}(\alpha)$  predicate that is true if  $i < c^n$  and  $\alpha$ , interpreted as a length- $c^n$  resolution proof for  $F$ , makes an invalid derivation in the  $i$ -th step. Note that this only depends on a constant number of nodes in the proof, and each node is described in  $\text{poly}(n)$  bits, hence  $\text{mistake}_F(n, \alpha, i)$  is indeed computable in deterministic polynomial time with oracle access to  $\alpha$ . The  $\forall\Sigma_1^b(\alpha)$ -sentence<sup>10</sup>

$$\forall n \in \text{Log} \exists i \leq c^n \text{mistake}_F(n, \alpha, i)$$

expresses the totality of the refuter problem as defined above; the provability of this sentence in relativized bounded arithmetic corresponds to the complexity of the refuter problem in  $\text{TFNP}^{\text{dt}}$ .<sup>11</sup>

**Width lower bounds for resolution.** In this paper, we also study the refuter problems corresponding to *width* lower bounds for resolution. Let  $F$  be a tautology without width- $w_F$  resolution proofs, the refuter problem for this width lower bound would be denoted as

$$\text{REFUTER}(w(F \vdash_{\text{Res}} \perp) \leq w_F).$$

The formalization of width lower bounds is essentially the same as that of size lower bounds, with the only difference that we now impose that every clause in the input resolution proof contains at most  $w_F$  literals. This can be done *syntactically* by only allocating  $w_F$  literals to each node.

### 1.2.2 Retraction Weak Pigeonhole Principles

This paper demonstrates that the complexity of refuter problems corresponding to resolution size lower bounds is tightly linked to the new complexity class  $\text{rwPHP}(\text{PLS})$ . Therefore, we need to introduce this class before describing our results.

Here, “ $\text{rwPHP}$ ” stands for the *retraction weak pigeonhole principle*:

For any two functions  $f : [N] \rightarrow [2N]$  and  $g : [2N] \rightarrow [N]$ , the function  $f \circ g : [2N] \rightarrow [2N]$  cannot be the identity function.

The term “retraction”, borrowed from category theory [Jeř07b], means that the principle concerns a pair of functions  $f, g$  where  $g$  is a “retraction”; the term “weak” indicates that the domain of  $g$  ( $[2N]$ ) is *much* larger than its range ( $[N]$ ). This principle, along with other variants of weak pigeonhole principles, is widely studied in the context of bounded arithmetic [PWW88, Kra01, MPW02, Tha02, Ats03, Kra04, Jeř04, Jeř07b, CLO24] and total search problems [KKMP21, Kor21, Kor22]; it is sometimes also called the “witnessing weak pigeonhole principle (WPHPWIT)” [Jeř07a, CLO24] and “LOSSY-CODE” [Kor22]. Clearly,  $\text{rwPHP}$  corresponds to a  $\text{TFNP}^{\text{dt}}$  problem: given (query access to) two functions  $f : [N] \rightarrow [2N]$  and  $g : [2N] \rightarrow [N]$ , find an input  $y \in [2N]$  such that  $f(g(y)) \neq y$ .

<sup>10</sup>Roughly speaking, a  $\forall\Sigma_1^b$ -sentence (resp.  $\forall\Sigma_1^b(\alpha)$ -sentence) is a sentence of the form

$$\forall x \exists y \varphi(x, y),$$

where  $|x|, |y|$  are polynomially related and  $\varphi$  is a polynomial-time relation (resp. polynomial-time relation with  $\alpha$  oracle); these sentences naturally express problems in  $\text{TFNP}$  (resp.  $\text{TFNP}^{\text{dt}}$ ). The notation  $n \in \text{Log}$  means that  $n$  is the length of some number, thus allowing one to reason about integers of magnitude  $2^{\text{poly}(n)}$  and strings of length  $\text{poly}(n)$ . In our particular case, it allows the length of the purported proof to be exponential in  $n$ . These are standard notations in bounded arithmetic.

<sup>11</sup>As a technical detail, we can also allow  $\alpha$  to take *parameters*  $\vec{z}$  that can be thought of as non-uniformity. That is, for each  $i \in [c^n]$ ,  $\alpha(\vec{z}, i)$  is the  $i$ -th bit of the proof. We consider the sentence

$$\forall n \in \text{Log}, \vec{z} \exists i \leq c^n \text{mistake}_F(n, \alpha(\vec{z}, \cdot), i)$$

which expresses that the proof encoded by  $\alpha(\vec{z}, \cdot)$  is not a valid length- $c^n$  resolution proof for  $F$ . The power of many natural principles with and without parameters are very different (see e.g., [ILW23, Section 4.3]).

Let  $\mathcal{P}$  be a problem in  $\text{TFNP}^{\text{dt}}$ , then one can define a class  $\text{rwPHP}(\mathcal{P})$  capturing the retraction weak pigeonhole principle where, informally speaking, the retraction function  $g$  can be computed in  $\mathcal{P}$ . In the decision tree model, the inputs of  $\text{rwPHP}(\mathcal{P})$  consist of:

1. (the evaluation table of) a function  $f : [N] \rightarrow [2N]$ , and
2.  $2N$  instances of  $\mathcal{P}$ , denoted as  $\{I_y\}_{y \in [2N]}$ , where each valid solution  $ans$  of each  $I_y$  is marked with an integer  $g_{y,ans} \in [N]$ .

The goal is to find an integer  $y \in [2N]$  along with a solution  $ans$  of  $I_y$  such that  $f(g_{y,ans}) \neq y$ . It is not hard to see that if  $\mathcal{P} \in \text{TFNP}^{\text{dt}}$  then  $\text{rwPHP}(\mathcal{P}) \in \text{TFNP}^{\text{dt}}$  (Fact 3.8). Furthermore,  $\text{rwPHP}(\mathcal{P})$  can be solved by a simple *randomized* algorithm given oracle access to any solver of  $\mathcal{P}$ .

The class  $\text{rwPHP}(\text{PLS})$  is defined as the problems reducible to  $\text{rwPHP}(\mathcal{P})$  for a  $\text{PLS}$ -complete problem  $\mathcal{P}$ . It can be shown that  $\text{rwPHP}(\text{PLS})$  does not depend on the exact choice of the  $\text{PLS}$ -complete problem  $\mathcal{P}$  (Fact 3.11).

**Witnessing for  $\text{T}_2^1 + \text{dwPHP}(\text{PV})$ .** Although  $\text{rwPHP}(\text{PLS})$  seems to be new to the  $\text{TFNP}$  community, it already appeared implicitly in the literature of bounded arithmetic. This class captures the  $\text{TFNP}$  problems whose totality is provable in  $\text{T}_2^1 + \text{dwPHP}(\text{PV})$ . In other words,  $\text{rwPHP}(\text{PLS})$  corresponds to the *witnessing theorem* for  $\text{T}_2^1 + \text{dwPHP}(\text{PV})$  (just like how  $\text{PLS}$  corresponds to a witnessing theorem for  $\text{T}_2^1$  [BK94]). This was noticed in [BKT14] where they showed every  $\forall\Sigma_1^b$ -consequence of  $\text{T}_2^1 + \text{dwPHP}(\text{PV})$  *randomly* reduces to  $\text{PLS}$ ; in fact, the same argument implies a deterministic reduction to  $\text{rwPHP}(\text{PLS})$ .

*Remark 2* (How Strong is  $\text{rwPHP}(\text{PLS})$ ?).

Since  $\text{rwPHP}(\text{PLS})$  can be seen as a randomized version of  $\text{PLS}$  (where the guarantee that “most randomness is good” is provided by the dual weak pigeonhole principle), its position in the  $\text{TFNP}^{\text{dt}}$  hierarchy is roughly the same as, but slightly higher than  $\text{PLS}$ . In particular, in the decision tree setting, it follows from the previous separations ( $\text{PLS} \not\subseteq \text{PPP}$  [GHJ<sup>+</sup>22] and  $\text{PLS} \not\subseteq \text{PPA}$  [BM04]) that  $\text{rwPHP}(\text{PLS})$  is contained in neither  $\text{PPP}$  nor  $\text{PPA}$ . Note that there is already a decision tree separation between  $\text{PLS}$  and the  $\text{TFNP}^{\text{dt}}$  problem corresponding to  $\text{rwPHP}$  (which follows from a resolution width lower bound for  $\text{rwPHP}$  [PT19, Proposition 3.4]), hence in the decision tree setting,  $\text{rwPHP}(\text{PLS})$  strictly contains  $\text{PLS}$ .

We also note that  $\text{T}_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha))$  is a relatively weak theory in the realm of relativized bounded arithmetic.<sup>a</sup> This theory is a subtheory of both  $\text{T}_2^2(\alpha)$  and Jeřábek’s (stronger) fragment for approximate counting  $\text{APC}_2(\alpha)$  [Jeř09]. It is also “weak” in the sense that unconditional unprovability results are known: it cannot prove the ordering principle [AT14] and the pigeonhole principle [PT19].

<sup>a</sup>The reader might have encountered claims in the literature that even weaker theories such as  $\text{S}_2^1$  or  $\text{APC}_1$  are “strong”, so it might be confusing for a reader unfamiliar with bounded arithmetic that we are claiming  $\text{T}_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha))$  as a “weak” theory. The reason is *relativization*: In our formalization, the purported resolution proof  $\alpha$  has *exponential* size, and we are only allowed to reason about objects in  $\text{PH}$  (think of  $\text{AC}^0$  circuits over  $\alpha$ ). This is much weaker than the setting where the proof  $\alpha$  has *polynomial* size and we are allowed to reason about polynomial-time concepts. This is roughly analogous to classifying the circuit class  $\text{AC}^0$  (i.e., relativized  $\text{PH}$ ) as “weak” and  $\text{P}/\text{poly}$  as “strong”.

### 1.3 Our Results

Our main results can be categorized into three parts: (1) bounded reverse mathematics ( $\text{TFNP}$  characterizations) for (several) resolution width lower bounds; (2) bounded reverse mathematics ( $\text{TFNP}$  characterizations) for (several) resolution size lower bounds; and (3) further applications in  $\text{TFNP}$  and proof complexity. We will describe the results related to width lower bounds first in Section 1.3.1, not only because they serve as prerequisites for the results regarding size lower bounds (discussed in Section 1.3.2), but also because the techniques therein find additional applications in  $\text{TFNP}$  and proof complexity (detailed in Section 1.3.3).

### 1.3.1 Refuters for Resolution Width Lower Bounds

The main message in this subsection is that the refuter problems corresponding to resolution width lower bounds are complete for the well-studied class PLS, the first syntactic subclass of TFNP introduced in the literature [JPY88].

We begin with the results related to the pigeonhole principle. The attentive reader may notice a subtle issue when formulating the refuter problem of width lower bound:  $\text{PHP}_{(n+1) \rightarrow n}$  already contains an axiom with width  $n$ , and the width lower bound for proving it is  $n$  as well. Thus, the corresponding width refuter problem becomes trivial. To address this, we instead consider the width refuter problem for a *constant-width analog* of  $\text{PHP}_{(n+1) \rightarrow n}$ , called  $\text{EPHP}_{(n+1) \rightarrow n}$ , which has constant-width axioms and an  $n/3$  width lower bound as shown in [BW01]. We characterize the complexity of its corresponding refuter problem:

**Theorem 5.2.**  $\text{REFUTER}(w(\text{EPHP} \vdash_{\text{Res}} \perp) < n/3)$  is PLS-complete.

A similar PLS-completeness result also holds for Tseitin formulas (on expander graphs), where  $e(G)$  below is the *expansion* parameter of the graph  $G$  (Definition 6.6).

**Theorem 6.9.**  $\text{REFUTER}(w(\text{Tseitin} \vdash_{\text{Res}} \perp) < e(G))$  is PLS-complete.

The techniques used in these results will be further extended to the refuter problems corresponding to black-box TFNP separations, specifically  $\text{PLS} \not\subseteq \text{PPP}$  and  $\text{PLS} \not\subseteq \text{PPA}$ , as described in Theorem 7.11 below.

To tackle Problem 1.1 though, we have to delve into the proofs of the exponential (size) lower bound. A *monotonized* version of the *width* lower bound plays a crucial role in the simplified proof by Beame and Pitassi [BP96]. In particular, they show that any resolution refutation of  $\text{PHP}_{(n+1) \rightarrow n}$  contains a clause  $C$  with “monotone width” of at least  $2n^2/9$  (see Section 4.1). We similarly characterize the complexity of its corresponding refuter problem (where the subscript *mono* denotes the monotone analog of the width refuter problem; the formal definition is provided in Section 4.1):

**Theorem 5.3.**  $\text{REFUTER}(w_{\text{mono}}(\text{PHP}_{(n+1) \rightarrow n} \vdash_{\text{Res}} \perp) < 2n^2/9)$  is PLS-complete.

This result serves as a key step toward addressing the size refuter problem for the pigeonhole principle, which will be discussed in the next subsection.

The PLS-hardness parts of all three results above stem from a *unified and simple proof*, detailed in Theorem 5.1. Conversely, the PLS-membership of these refuter problems is established by carefully analyzing the proofs in [BP96, BW01] and demonstrating that “PLS-reasoning” suffices to prove these lower bounds. (In fact, these proofs can be formalized in the theory  $\text{T}_2^1(\alpha)$ , and the PLS-membership follows directly from the witnessing theorem in [BK94].)

**A non-uniform universal PLS-membership.** Finally, we establish a *universal* PLS-membership result with respect to *non-uniform* decision tree reductions: for any resolution width lower bound against every unsatisfiable CNF, *as long as the lower bound is correct*, the corresponding refuter problem can be reduced to PLS under *non-uniform* decision tree reductions.

**Theorem 6.1.** *Let  $\mathcal{F}$  be any (possibly non-uniform) family of unsatisfiable CNFs with polynomially many clauses, and let  $w_0$  be any valid resolution width lower bound for  $\mathcal{F}$ . Then there exists a (non-uniform) decision-tree reduction from  $\text{REFUTER}(w(\mathcal{F} \vdash_{\text{Res}} \perp) < w_0)$  to PLS.*

Both the formulation and proof of this result inherently require non-uniformity for at least two reasons: (1) it is computationally hard to check whether an arbitrarily given CNF is unsatisfiable, and (2) even assuming that the given CNF is unsatisfiable, it is hard to calculate the resolution width lower bound. See Section 6.1 for further discussion.

*Remark 3* (Uniform vs. non-uniform reductions). Note that if one only cares about non-uniform reductions, then (the PLS-membership parts of) [Theorem 5.2](#) and [Theorem 6.9](#) are merely special cases of [Theorem 6.1](#). Nevertheless, we believe that the uniform PLS-membership results in [Theorem 5.2](#) and [Theorem 6.9](#) are informative, as they actually show that the corresponding lower bounds can be formalized in  $T_2^1(\alpha)$ ; in fact, the *code* of the Turing machine that implementing the uniform reduction to PLS effectively acts as a *proof* of the width lower bound using a *local search* argument. They are also crucial for the uniform rwPHP(PLS)-memberships for the size refuter problems. However, the decision tree reduction in [Theorem 6.1](#) seems to require  $\exp(n)$  bits of non-uniformity, making it *highly* non-uniform.

On the other hand, the non-uniform reduction in [Theorem 6.1](#) implies an intriguing proof complexity upper bound: *Small-width resolution can prove width lower bounds for resolution itself!* (See [Section 1.3.3](#) for more details.) Uniformity is not required for this application, allowing us to derive more proof complexity upper bounds using [Theorem 6.1](#): *every* resolution width lower bound *that is correct* can be proved in low-width resolution. (The size lower bound analog of [Theorem 6.1](#) remains unknown, hence we can only show proof complexity upper bounds for tautologies encoding *specific* resolution size lower bounds.)

### 1.3.2 Refuters for Resolution Size Lower Bounds

Our main message in this subsection is that the refuter problems corresponding to many resolution size lower bounds are complete for rwPHP(PLS), the TFNP subclass introduced in [Section 1.2.2](#). Indeed, the theorems presented in this subsection suggest that rwPHP(PLS) captures the complexity of proving *the easiest-to-prove* size lower bounds for resolution. Our workflow is the same as before:

- First, we show that for many notable resolution size lower bounds proven in the literature, the corresponding refuter problems reduce to rwPHP(PLS). Specifically, we identify a common technique for proving resolution size lower bounds, which we call “random restriction + width lower bounds”, and demonstrate that if a resolution size lower bound can be proven using it, then the corresponding refuter problem generally falls within rwPHP(PLS).
- Next, we present a *unified* rwPHP(PLS)-hardness result: the refuter problems for resolution size lower bounds are rwPHP(PLS)-hard, and the hardness proof *does not* depend on the hard tautology considered. Thus, we conclude the rwPHP(PLS)-completeness of many refuter problems for resolution size lower bounds.

The rwPHP(PLS)-hardness of size lower bound refuters turns out to be more challenging than the PLS-hardness of width lower bound refuters, as discussed in [Section 2.2](#).

We begin by showing that [Problem 1.1](#) reduces to rwPHP(PLS):

**Theorem 1.4** (Informal version of [Theorem 4.16](#)). *There exists an absolute constant  $c > 1$  and an efficient decision-tree reduction from the problem  $\text{REFUTER}(s(\text{PHP}_{(n+1) \rightarrow n} \vdash_{\text{Res}} \perp) \leq c^n)$  to rwPHP(PLS).*

In fact, we show that  $T_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha))$  proves the sentence

$$\forall n \in \text{Log} \exists i \leq c^n \text{mistake}_{\text{PHP}}(n, \alpha, i),$$

i.e.,  $\alpha$  is not a length- $c^n$  resolution proof for PHP, by formalizing the classical proofs in [[Hak85](#), [CP90](#), [BP96](#)]; [Theorem 4.16](#) then follows from the witnessing theorem for  $T_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha))$ . In the technical overview ([Section 2.1](#)) and the main proof ([Section 4.2](#)), we present the reduction from the refuter problem  $\text{REFUTER}(s(\text{PHP} \vdash_{\text{Res}} \perp) \leq c^n)$  to rwPHP(PLS) directly, without relying on witnessing theorems.

It turns out that a large variety of resolution size lower bounds can be proven using the paradigm of “random restriction + width lower bounds,” including those for XOR-lifted formulas [[DR03](#)], Tseitin formulas [[Urq87](#), [Sch97](#)], and random CNFs [[CS88](#)]. We show that all these lower bounds have corresponding refuter problems in rwPHP(PLS) (see [Theorem 6.4](#), [Theorem 6.12](#), and [Theorem 6.13](#), respectively). These

results provide strong evidence that  $\text{rwPHP(PLS)}$  (or  $\text{T}_2^1(\alpha) + \text{dwPHP(PV}(\alpha))$ ) captures the “complexity” of this popular proof technique for resolution lower bounds.

We complement the above results by showing that for *every* unsatisfiable family of CNFs  $\{F_n\}$  that requires resolution size greater than  $s_F(n)$ , the corresponding refuter problem  $\text{REFUTER}(s(F_n \vdash_{\text{Res}} \perp) \leq s_F(n))$  is hard for  $\text{rwPHP(PLS)}$ .

**Theorem 1.5** (Informal version of [Theorem 5.4](#)). *For every unsatisfiable family of CNF formulas  $\{F_n\}$  and parameter  $s_F(n)$  such that every resolution refutation of  $F_n$  requires more than  $s_F(n)$  clauses, there exists a decision tree reduction of depth  $\text{poly}(n)$  from  $\text{rwPHP(PLS)}$  to  $\text{REFUTER}(s(F_n \vdash_{\text{Res}} \perp) \leq s_F(n))$ .<sup>12</sup>*

Note that [Theorem 1.5](#) holds for *every* hard tautology, whereas the  $\text{rwPHP(PLS)}$  upper bounds such as [Theorem 1.4](#) are only known to hold for some natural families of hard tautologies. For these natural tautologies, we establish a *reversal* in the bounded reverse mathematics of proof complexity lower bounds: The power of “ $\text{rwPHP(PLS)}$ -reasoning” is *sufficient* for implementing a popular proof strategy that can prove all these resolution lower bounds and, at the same time, is *necessary* for proving *any* resolution lower bound.

*Remark 4.* We also note that [Theorem 1.5](#) requires decision tree depth  $\text{poly}(n)$  regardless of  $s_F$ , and is thus only considered “efficient” when  $s_F = 2^{n^{\Omega(1)}}$ . However, this is merely an artifact of our definition of “efficiency” in the decision tree setting, i.e., if the input length is  $N$ , then depth- $\text{polylog}(N)$  decision trees are considered “efficient”. In fact, even if  $s_F = 2^{n^{o(1)}}$ , each node in the purported length- $s_F$  resolution proof still requires  $\text{poly}(n)$  bits to represent, so it takes  $\text{poly}(n)$  query complexity to *verify* a solution of the refuter problem. Therefore, it still makes sense *in the particular setting of refuter problems* to consider a decision tree reduction *efficient* if its query complexity is at most  $\text{poly}(n)$ . We interpret [Theorem 1.5](#) to mean that “ $\text{rwPHP(PLS)}$ -reasoning” is necessary for proving *not only subexponential but any moderately large size* lower bound for resolution.

The proof of [Theorem 1.5](#) is heavily inspired by the NP-hardness of automating resolution [[AM20](#)] and the exposition of this result in [[dRGN<sup>+</sup>21](#)]. In these proofs, it was crucial to show that resolution cannot prove lower bounds against itself; in particular, [[dRGN<sup>+</sup>21](#), Section 5] showed that resolution requires a large (block-)width to prove resolution lower bounds. Notably, the proof in [[dRGN<sup>+</sup>21](#)] is by a reduction from  $\text{rwPHP}$ , i.e., resolution cannot prove lower bounds against itself because resolution cannot prove  $\text{rwPHP}$ . We strengthen these results by reducing a stronger problem — $\text{rwPHP(PLS)}$  instead of  $\text{rwPHP}$  — to the refuter problems, thereby obtaining a *tight* characterization of these refuter problems.

Finally, our results provide an intriguing characterization of the provably total NP search problems in  $\text{T}_2^1 + \text{dwPHP(PV)}$  (see [Corollary 5.7](#)). That is:

Just as “every DAG has a sink” characterizes the  $\forall \Sigma_1^b$ -consequences of  $\text{T}_2^1$  [[BK94](#)], “resolution requires  $2^{\Omega(n)}$  size to prove PHP” characterizes the  $\forall \Sigma_1^b$ -consequences of  $\text{T}_2^1 + \text{dwPHP(PV)}$ .

### 1.3.3 Applications

Besides being interesting in itself, our study of refuter problems also reveals several new insights into these well-studied proof complexity lower bounds and TFNP separations. More specifically, we translate our results into different languages using the generic connection between  $\text{TFNP}^{\text{dt}}$  and proof complexity via the *false clause search* problem (see, e.g., [[dRGR22](#)]): For an unsatisfiable CNF  $F = C_1 \wedge \dots \wedge C_M$ , the false clause search problem  $\text{Search}(F)$  is a  $\text{TFNP}^{\text{dt}}$  problem where, given oracle access to an input  $x \in \{0, 1\}^N$ , the goal is to find a clause  $C_i$  such that  $C_i(x) = \text{false}$ . Any  $\text{TFNP}^{\text{dt}}$  problem can be written as a false clause search problem for a family of low-width CNFs, and vice versa. In particular, a family of unsatisfiable CNFs has low-width resolution refutations if and only if the corresponding false clause search problem reduces to PLS [[Raz95b](#)] (see also [[Kam19](#), Section 8.2.2] for an exposition). See [Figure 1](#) for a diagram that summarizes the translations of our main results in different languages.

<sup>12</sup>This theorem requires a mild technical condition that  $s_F(n)$  should be moderately larger than the size of the  $\text{rwPHP(PLS)}$  instance; see the formal statement in [Theorem 5.4](#) for details.

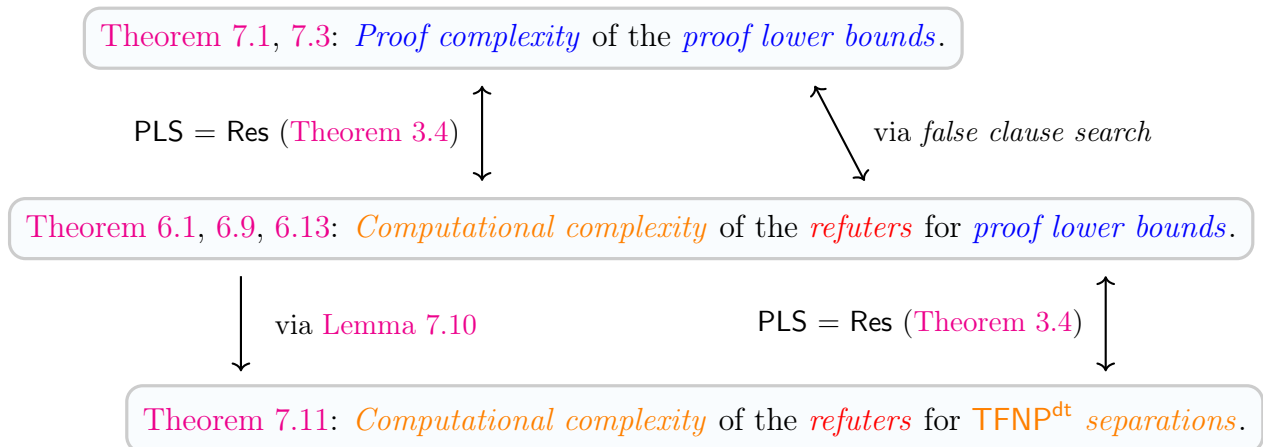


Figure 1: Translations of the main results in different languages. A one-way arrow represents an implication, and a two-way arrow indicates an equivalence.

**Proof complexity of proof lower bounds.** We first use our results to provide surprisingly efficient proofs for proof complexity lower bounds. Note that a proof complexity lower bound can be expressed by a family of CNFs  $\mathcal{F}_{\text{LB}}$  by formulating the corresponding refuter problem as a false clause search problem  $\text{Search}(\mathcal{F}_{\text{LB}})$  (see, e.g., Section 7.1).

In particular, there exists a family of  $\tilde{O}(w)$ -width CNFs that encodes a width- $w$  resolution lower bound. Then, since PLS and low-width resolution are equivalent, Corollary 6.2 implies the following *upper bound* on the resolution width required to prove resolution lower bounds.

**Theorem 1.6** (Informal version of Theorem 7.1). *Any width- $w$  resolution lower bound can be proved in resolution width  $\tilde{O}(w)$ .*

We also use our rwPHP(PLS) upper bounds to show that poly( $n$ )-width *random resolution* [BKT14, PT19] can prove exponential-size resolution lower bounds (encoded as poly( $n$ )-width CNFs). In fact, using our results on random  $k$ -CNFs (Theorem 6.13), we can show that *most* resolution size lower bounds are provable in low-width random resolution:

**Theorem 1.7** (Informal version of Theorem 7.3). *With high probability over a random  $k$ -CNF  $F$ , the resolution size lower bound  $s(F \vdash_{\text{Res}} \perp) > 2^{\Omega(n)}$  can be proved in random resolution width of poly( $n$ ).*

These results stand in stark contrast with Garlík’s result [Gar19] that tautologies encoding any resolution size lower bounds are hard for resolution: We show that either switching to width lower bounds (Theorem 1.6) or considering *random resolution* (Theorem 1.7) makes these lower bound tautologies easy to prove!<sup>13</sup>

**Complexity of refuting black-box TFNP separations.** We also consider the refuter problem for *black-box TFNP separations*. Let  $A, B$  be two TFNP<sup>dt</sup> classes such that  $A \not\subseteq B$ . Informally,  $\text{Ref}(A \subseteq B)$  is the class of problems reducible to the following kind of “refuter” problems: The input is a purported decision tree reduction from  $A$  to  $B$ , and the solution is a short witness showing that the reduction is wrong. The refuter problems for TFNP<sup>dt</sup> separations also lie in TFNP<sup>dt</sup>, as their totality follows from the correctness of the black-box separation  $A \not\subseteq B$ .

<sup>13</sup>In our Section 7.1, the lower bound tautologies use *binary encoding*, where (e.g.) the predecessors of every node are encoded by  $O(\log N)$  bits. In contrast, Garlík [Gar19] uses *unary encoding* where for every pair of nodes  $(i, j)$  (a minor detail is that [Gar19] requires  $i$  to be “one level above”  $j$ ), there is a Boolean variable  $x_{i,j}$  indicating whether  $i$  is a predecessor of  $j$ . As [Gar19] pointed out, a resolution lower bound for the unary-encoded refutation statements implies a similar lower bound for the binary-encoded refutation statements. On the other hand, since we are proving *width* upper bounds and the unary encoding already results in large-width CNFs, we can only afford to use binary encoding (see Remark 6).

The complexity of such refuter problems measures the strength of the arguments used for black-box separation results. For example, the following corollary conveys a simple but often overlooked fact: *when separating a syntactic  $\text{TFNP}^{\text{dt}}$  subclass  $A$  from  $B$ , it is necessary to incur the totality principle of  $B$ .*

**Corollary 7.9.** *For any two  $\text{TFNP}^{\text{dt}}$  classes  $A, B$  such that  $A \not\subseteq B$ ,  $B \subseteq \text{Ref}(A \subseteq B)$ .*

Due to the connection between  $\text{TFNP}^{\text{dt}}$  and proof complexity, the refuter problem for each black-box  $\text{TFNP}$  separation naturally aligns with a corresponding refuter problem for a proof complexity lower bound. In particular, we build a uniform reduction from the refuter problems for separations from PLS to the refuter problems for resolution width lower bounds (Lemma 7.10), because showing a  $\text{TFNP}^{\text{dt}}$  subclass  $A$  is not in PLS is essentially showing a resolution width lower bound for the formula expressing the totality of  $A$ .

Note that the false clause search problem for EPHP and Tseitin are in PPP and PPA respectively. Therefore, using our characterization of the resolution width refuter for EPHP (Theorem 4.3) and Tseitin (Theorem 6.9), we conclude that *it is necessary and sufficient to use local search principle to separate PPP and PPA from PLS in the black-box setting.*

**Theorem 1.8** (Informal version of Theorem 7.11).  $\text{Ref}(\text{PPP} \subseteq \text{PLS}) = \text{Ref}(\text{PPA} \subseteq \text{PLS}) = \text{PLS}$ .

## 1.4 Discussions, Speculations, and Future Directions

This paper initiates a research program that attempts to understand, for every proof system  $\mathcal{P}$  of interest, the metamathematics of proving lower bounds against  $\mathcal{P}$  through the lens of refuter problems.<sup>14</sup> Our results on resolution suggest that this is a promising direction. There are a plethora of future research directions, both regarding “weak” systems (where we already know strong lower bounds against  $\mathcal{P}$ -proofs) and “strong” ones (where we are still struggling to prove non-trivial lower bounds against  $\mathcal{P}$ ).

**Weak proof systems.** It might be feasible to characterize the complexity of refuter problems for weak proof systems. How does the complexity of refuting lower bounds for  $\mathcal{P}$  compare with  $\mathcal{P}$  itself (or, more precisely, the  $\text{TFNP}^{\text{dt}}$  subclass corresponding to  $\mathcal{P}$  [BFI23])? In the case that  $\mathcal{P}$  is resolution, our work shows that the complexity of refuting *width* lower bounds for  $\mathcal{P}$  is exactly  $\mathcal{P}$  itself (i.e., PLS), and the complexity of refuting *size* lower bounds is a randomized version of  $\mathcal{P}$  (i.e.,  $\text{rwPHP}(\text{PLS})$ ). Thus, it seems reasonable to conjecture that for “weak” proof systems, the complexity of proving lower bounds against them is not much higher than themselves.

Moreover, the proof complexity of proof complexity lower bounds is intimately connected to the hardness of automatability of proof systems, see e.g., [AM20, GKMP20, Bel20, dRGN<sup>+</sup>21, IR22, Gar24, Pap24]. We expect that a thorough understanding of the former would help make progress on the latter as well.

**Strong proof systems.** The situation for strong proof systems seems much more mysterious. For strong proof systems  $\mathcal{P}$  (think of  $\mathcal{P}$  being Frege or Extended Frege), it is even unclear whether there should be an “easiest-to-prove” lower bound for  $\mathcal{P}$  (which would correspond to a syntactic subclass  $\mathcal{C}(\mathcal{P}) \subseteq \text{TFNP}^{\text{dt}}$  that characterizes the complexity of proving lower bounds for  $\mathcal{P}$ ). Even if such a  $\mathcal{C}(\mathcal{P})$  exists, it is unclear if it is captured within our current landscape of  $\text{TFNP}^{\text{dt}}$ .<sup>15</sup>

This suggests the following possibility: The reason that we have not been able to prove lower bounds for  $\mathcal{P}$  is that  $\mathcal{C}(\mathcal{P})$  is a very complicated class, far beyond our current understanding of  $\text{TFNP}^{\text{dt}}$  and

<sup>14</sup>We believe the similar research program for circuit lower bounds would also be fruitful, which has already started since [CJSW24, Kor22, CLO24] if not earlier. We limit our discussions to proof lower bounds here.

<sup>15</sup>Note that the question of where  $\mathcal{C}(\mathcal{P})$  sits in the  $\text{TFNP}^{\text{dt}}$  hierarchy is merely a restatement of the open problem of determining the proof complexity of proof complexity lower bounds for  $\mathcal{P}$ . For example,  $\mathcal{C}(\mathcal{P})$  is a subclass of  $\text{PTFNP}$  [GP18a] if and only if  $\text{Q-EFF}$  (the proof system underlying the definition of  $\text{PTFNP}$ ) can prove lower bounds for  $\mathcal{P}$ .

bounded arithmetic. An even more speculative hypothesis would be that the proof systems  $\mathcal{P}$  for which we are able to prove lower bounds are exactly those where  $\mathcal{C}(\mathcal{P})$  is not “much” higher than  $\mathcal{P}$  themselves. We hope that future work will determine to what extent these hypotheses are correct.

The case of  $\text{AC}^0[p]$ -Frege (where  $p$  is a prime) is of particular interest. Although strong lower bounds for  $\text{AC}^0[p]$  circuits have been known for decades [Raz87, Smo87], we have not yet succeeded in turning these circuit lower bounds into proof complexity lower bounds against  $\text{AC}^0[p]$ -Frege (see, e.g., [MP96, BKZ15]). The paper [BIK<sup>+</sup>97] laid out a research program towards  $\text{AC}^0[p]$ -Frege lower bounds by studying weaker algebraic proof systems such as the Nullstellensatz [BIK<sup>+</sup>94] and Polynomial Calculus [CEI96, Raz98]. After a few decades, we have become proficient at proving lower bounds against such algebraic proof systems, but lower bounds against  $\text{AC}^0[p]$ -Frege remain elusive. Is it because the refuter problems corresponding to  $\text{AC}^0[p]$ -Frege lower bounds are *fundamentally different* from those for the weaker algebraic proof systems? Does our metamathematical  $\text{TFNP}^{\text{dt}}$  perspective bring new insights to this long-standing open question?

## 1.5 Further Related Works

**Refuter problems for circuit lower bounds.** Our study of the refuter problems for proof lower bounds is strongly influenced by the line of work on refuter problems for circuit lower bounds. Chen, Jin, Santhanam, and Williams [CJSW24] call a lower bound *constructive* if the corresponding refuter problem can be solved in deterministic polynomial time, and they argued that constructivity is a desirable aspect of lower bounds. Chen, Tell, and Williams [CTW23] showed that for many lower bounds against randomized computational models, their refuter problems characterize derandomizing  $\text{pr-BPP}$ . The main result of Korten [Kor22] can also be seen as the  $\text{WPHPWIT}$ -hardness of refuter problems for one-tape Turing machine lower bounds. Pich and Santhanam [PS26] showed how to turn proof complexity lower bounds into circuit lower bounds, assuming the refuter problem for the (conjectured) lower bound  $\text{SAT} \notin \text{P/poly}$  is “provably easy” in a certain sense. Finally, the results of Chen, Li, and Oliveira [CLO24] can be interpreted as the  $\text{PWPP}$ - and  $\text{WPHPWIT}$ -completeness of various refuter problems.

It is also worth mentioning that Ebtehaj [Ebt23] studied the refuter problems for  $\mathcal{A} \not\subseteq \text{BPP}$  for each (type-1) subclass  $\mathcal{A} \subseteq \text{TFNP}$  that is indeed hard. However, [Ebt23] did not obtain any completeness results for such refuter problems.

**Unprovability of complexity upper bounds.** In parallel to the investigation of unprovability of complexity lower bounds, there is another line of work showing the unprovability of complexity *upper* bounds in fragments of bounded arithmetic [CK07, KO17, BKO20, BM20, CKKO21, ABM23]. For example, Krajíček and Oliveira [KO17] proved that Cook’s theory  $\text{PV}$  cannot prove  $\text{P} \subseteq \text{SIZE}[n^k]$ , and Atserias, Buss, and Müller [ABM23] proved that the theory  $\text{V}_2^0$  cannot prove  $\text{NEXP} \subseteq \text{P/poly}$ . These results are equivalent to the *consistency* of lower bounds with fragments of bounded arithmetic, thus in some sense representing progress towards proving circuit lower bounds.<sup>16</sup> Indeed, [CKKO21] presented a general framework for showing such consistency results by proving lower bounds against circuits with a certain uniformity condition called “ $\text{LEARN}$ -uniformity”, and the techniques employed in many of these papers are inspired by uniform circuit lower bounds such as [SW14].

**Witnessing theorems.**  $\text{TFNP}$  and bounded arithmetic are connected through *witnessing theorems*: each theory is associated with the class of  $\text{TFNP}$  problems whose totality is provable in this theory.

<sup>16</sup>The “conventional wisdom” seems to believe that the complexity lower bounds are true (for discussions, see <https://rjlipton.com/conventional-wisdom-and-ppp/>, accessed Mar 14, 2026). Hence, unprovability of complexity lower bounds can be seen as the difficulty for proving this “conventional wisdom”, while unprovability of complexity upper bounds represents progress towards proving it. One should keep in mind that the opposite opinion makes equal sense: for a believer of complexity *upper* bounds, the unprovability of these upper bounds indicates the difficulty of confirming their belief, while the unprovability of lower bounds implies progress towards it!

Perhaps the best-known witnessing theorem is Buss’s one [Bus85]: every NP search problem provably total in  $S_2^1$  can be solved in deterministic polynomial time. The class PLS and its generalizations such as CPLS capture the NP search problems provably total in higher levels of bounded arithmetic hierarchy [BK94, KST07, ST11, PT12]; in this sense, witnessing theorems also provide a systematic method for defining new syntactic subclasses of TFNP. Other witnessing theorems considered in the literature include [KNT11, BB17, KT22]. Our paper contributes to this line of research by characterizing the class of NP search problems provably total in  $T_2^1 + \text{dwPHP}(\text{PV})$  by the refuter problems corresponding to many resolution lower bounds, in particular the problem  $\text{REFUTER}(s(\text{PHP}_{(n+1) \rightarrow n} \vdash_{\text{Res}} \perp)) < c^n$ .

**Comparison with the consistency search problem.** We note that the refuter problem looks superficially similar to  $\text{WRONGPROOF}$ , the *consistency search* problem for proof systems [BB17, GP18a, Pud20]. Let  $\mathcal{P}$  be a proof system,  $\text{WRONGPROOF}(\mathcal{P})$  is the  $\text{TFNP}^{\text{dt}}$  problem that given as input a purported  $\mathcal{P}$ -proof  $\Pi$  of an *incorrect statement*, asks for the location of an invalid derivation in  $\Pi$ .

Although both  $\text{WRONGPROOF}$  and our refuter problems take a purported proof as input and ask for an invalid derivation in the proof, we think that these two problems are fundamentally different, because they have different *reasons of totality*. Roughly speaking, the totality of  $\text{WRONGPROOF}$  is proved by the *soundness* of  $\mathcal{P}$ , and the totality of  $\text{REFUTER}$  is guaranteed by *lower bound proofs*. We elaborate on this in [Appendix B](#).

Another (superficial) similarity between these two problems is that both problems are used to characterize the provably total NP problems in bounded arithmetic. The consistency search problems for Frege and Extended Frege characterize the  $\forall\Sigma_1^b$ -consequences of  $U_2^1$  and  $V_2^1$  respectively [BB17], while in this paper we show that the refuter problem for resolution (with a suitable hard tautology) characterizes the  $\forall\Sigma_1^b$ -consequences of  $T_2^1 + \text{dwPHP}(\text{PV})$ .

## Paper Organization

The main body of the paper is structured so that the initial sections primarily focus on presenting the complexities of refuter problems of the pigeonhole principle. Specifically, [Section 3](#) contains the necessary preliminaries and formal definitions. [Section 4](#) presents the complexity upper bound of the refuters for  $\text{PHP}_{(n+1) \rightarrow n}$ . [Section 5](#) complements the previous section by showing universal hardness results for refuting any narrow or short resolution proofs. In [Section 6](#), we extend our results to many other formulas, including XOR-lifted formulas, Tseitin formulas, and random  $k$ -CNFs. Finally, [Section 7](#) introduces two novel applications of our results in understanding the proof complexity of proof complexity lower bounds and the complexity of black-box TFNP separations.

## 2 Technical Overview

### 2.1 Refuter Problems in $\text{rwPHP}(\text{PLS})$

In this subsection, we explain how the lower bound proof in [CP90, BP96] yields a reduction from the problem  $\text{REFUTER}(s(\text{PHP}_{(n+1) \rightarrow n} \vdash_{\text{Res}} \perp) \leq c^n)$  to  $\text{rwPHP}(\text{PLS})$ . As mentioned before, this is essentially a formalization of the lower bound proof in  $T_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha))$ , and the reduction follows from the witnessing theorem for this theory. However, this subsection will describe the reduction without invoking the witnessing theorem (nor does the formal proof in [Section 4.2](#) use the witnessing theorem). We hope that by opening up the black box of the witnessing theorem, it would become clearer how each component in the proof corresponds to a component in the reduction to  $\text{rwPHP}(\text{PLS})$ .

The proof of [CP90, BP96] consists of two components:

- (Random restrictions) First, we carefully design a distribution of random restrictions  $\mathcal{R}$  under which the following holds. (1) With high probability over  $\rho \leftarrow \mathcal{R}$ , any fixed size- $c^n$  resolution proof will

simplify to a resolution proof of *width*<sup>17</sup> at most  $w$  under  $\rho$  (for some parameter  $w$ ); (2) the pigeonhole principle  $\text{PHP}_{(n+1) \rightarrow n}$  remains to be the pigeonhole principle (of a slightly smaller size  $\text{PHP}_{(n'+1) \rightarrow n'}$ ) under any restriction  $\rho \in \mathcal{R}$ . Moreover,  $\mathcal{R}$  is the uniform distribution over some set of restrictions; we abuse notation and use  $\mathcal{R}$  to also denote this set.

- (Width lower bound) Then, we invoke the width lower bound for the pigeonhole principle and show that resolution cannot prove  $\text{PHP}_{(n'+1) \rightarrow n'}$  in width  $w$ .

Given a resolution proof  $\Pi$  of size  $c^n$ , the fact that most restrictions simplify  $\Pi$  into a small-width proof can be shown by a *compression argument*: given a clause  $C_i \in \Pi$  and a restriction  $\rho \in \mathcal{R}$  that does *not* shrink  $C_i$  into a clause of width  $\leq w$ , one can describe  $\rho$  in  $\ell_{\text{comp}}$  bits for some small  $\ell_{\text{comp}}$ . In what follows, it suffices if  $\ell_{\text{comp}} \leq \log |\mathcal{R}| - \log(c^n) - 1$ , which is indeed the case under a suitable choice of parameters. Note that for comparison, if such a clause  $C_i$  were not known, it would require, information-theoretically, at least  $\log |\mathcal{R}|$  bits to encode any restriction  $\rho \in \mathcal{R}$ . This compression argument implies that a random  $\rho \leftarrow \mathcal{R}$  shrinks a fixed clause w.p.  $\geq 1 - \frac{1}{2c^n}$ , thus the existence of a good  $\rho$  shrinking the whole proof  $\Pi$  follows from a union bound over the  $c^n$  clauses in  $\Pi$ .

For every restriction  $\rho$ , one can compute a proof  $\Pi|_\rho$  with each clause  $C \in \Pi$  replaced by  $C|_\rho$ , the restriction of  $C$  under  $\rho$ ; if  $\text{width}(C|_\rho) > w$ , we truncate  $C|_\rho$  to force its width to be at most  $w$ . Since  $\Pi|_\rho$  is a width- $w$  resolution proof, it follows from the width lower bound that it does not prove  $\text{PHP}_{(n'+1) \rightarrow n'}$ .

Our reduction from  $\text{REFUTER}(s(\text{PHP}_{(n+1) \rightarrow n} \vdash_{\text{Res}} \perp)) \leq c^n$  to  $\text{rwPHP}(\text{PLS})$  works as follows.

- Let  $N := |\mathcal{R}|/2$ . The function  $f : [N] \rightarrow [2N]$  takes as inputs  $(i, s)$  where  $i \in [c^n]$  denotes a node in  $\Pi$  and  $s$  is the compressed description of a random restriction  $\rho$  that fails to simplify  $C_i$  to width  $w$  (note that this takes  $\log(c^n) + \ell_{\text{comp}} \leq \log N$  bits), and outputs the standard encoding of  $\rho$  (in  $\log |\mathcal{R}| = \log(2N)$  bits). It is easy to see that every restriction  $\rho$  outside the range of  $f$  would be a good restriction (that successfully shrinks every clause in  $\Pi$  into width  $w$ ).
- For each  $\rho \in \mathcal{R} \cong [2N]$ ,  $\Pi|_\rho$  is a width- $w$  resolution proof. By [Theorem 4.8](#), we can reduce the problem of finding an illegal derivation in  $\Pi|_\rho$  to PLS. Call this instance  $I_\rho$ .
- Finally, let  $i \in [c^n]$  be an illegal derivation in  $\Pi|_\rho$  (that can be found in PLS). There are two reasons that the  $i$ -th step is illegal in  $\Pi|_\rho$ : first, it might already be an invalid derivation in  $\Pi$ ; second, the width of  $C_i|_\rho$  might be greater than  $w$ , thus the error happens when we truncate  $C_i|_\rho$  to width  $w$ . In the second case, let  $s$  be the  $\ell_{\text{comp}}$ -bit description of  $\rho$  given that it does not simplify  $C_i$  to width  $\leq w$ , and  $g_{\rho,i} := (i, s)$ , then  $f(g_{\rho,i}) = \rho$ .

It follows that once we found any  $\rho$  and  $i$  such that  $i$  is an answer for  $I_\rho$  and  $f(g_{\rho,i}) \neq \rho$ , then the  $i$ -th step is illegal for the first reason stated above, i.e., the  $i$ -th step in  $\Pi$  is also invalid. See [Figure 2](#) for a high-level overview of this construction.

**Random restrictions + width lower bounds.** It turns out that the above proof template that combines random restrictions and width lower bounds is very popular in proving resolution lower bounds. Given a hard tautology  $F$ , we design a family of restrictions  $\mathcal{R}$  such that (1) Any fixed *short* resolution proof will simplify to a *narrow* resolution proof under  $\mathcal{R}$ , and (2) even after a random restriction in  $\mathcal{R}$ ,  $F$  remains hard for narrow resolution proofs. Note that the family  $\mathcal{R}$  is usually carefully chosen according to the hard tautology  $F$ ; e.g.,  $\mathcal{R}$  corresponds to *partial matchings* when  $F = \text{PHP}$  [[BP96](#)] and corresponds to *random edge sets* when  $F = \text{Tseitin}$  [[Sch97](#)].

As mentioned before, this proof strategy is capable of proving resolution size lower bounds for various hard tautologies, and we can use a similar argument as the above to show that the refuter problems

<sup>17</sup>In fact, in the case of PHP, we obtain a resolution proof of small *monotone width*. We omit the distinction between width and monotone width in the overview and refer the reader to [Section 4](#) for details.

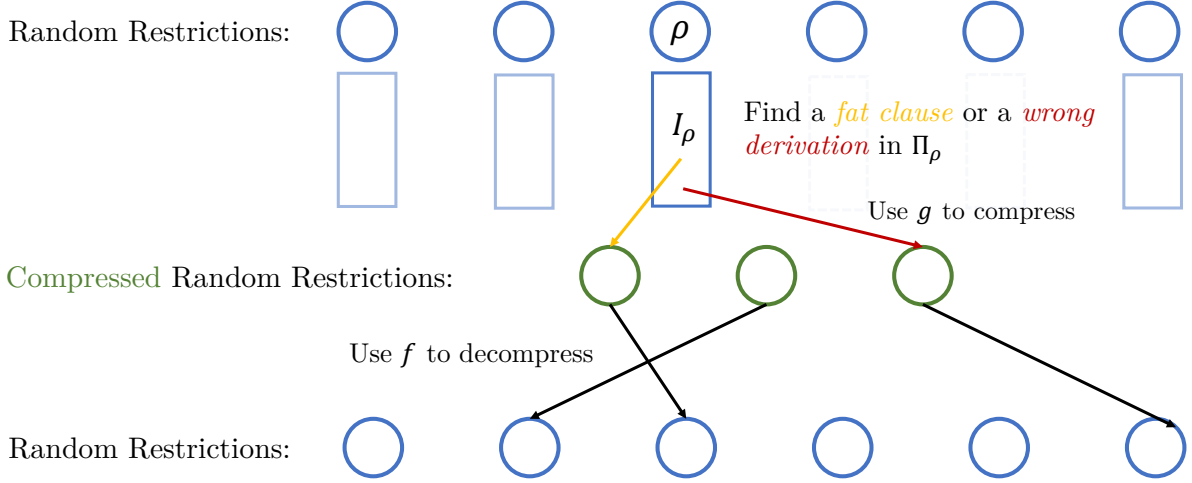


Figure 2: The  $\text{rwPHP(PLS)}$  instance constructed from  $\text{REFUTER}(s(\text{PHP}_{(n+1) \rightarrow n} \vdash_{\text{Res}} \perp) \leq c^n)$ .

corresponding to these resolution size lower bounds are in  $\text{rwPHP(PLS)}$ . This includes XOR-lifted formulas (Section 6.2), Tseitin tautologies (Section 6.3), and random  $k$ -CNFs (Section 6.4).

In fact, it is quite intuitive to formalize “random restrictions + width lower bounds” in  $\text{T}_2^1(\alpha) + \text{dwPHP(PV}(\alpha))$ . Roughly speaking, we first use  $\text{dwPHP(PV}(\alpha))$  to formalize the compression argument and show that most random restrictions will shrink the resolution proof (represented by  $\alpha$ ) into a narrow one; then we use  $\Sigma_1^b(\alpha)$ -MIN (which is available in  $\text{T}_2^1(\alpha)$ ) to prove a resolution width lower bound.

## 2.2 Refuter Problems are $\text{rwPHP(PLS)}$ -Hard

In this subsection, we explain the ideas behind the reduction from  $\text{rwPHP(PLS)}$  to the refuter problems for resolution size lower bounds. In fact, a reduction from  $\text{rwPHP}$  to the refuter problems already appeared in [dRGN<sup>+</sup>21], which provides a streamlined proof of the celebrated NP-hardness of automating resolution [AM20]. It turns out that with minor modifications, the same proof can be adapted to reduce not only  $\text{rwPHP}$  but also  $\text{rwPHP(PLS)}$  to the refuter problems, thereby proving Theorem 5.4. Hence, the remainder of this subsection will focus on the  $\text{rwPHP}$ -hardness result from [dRGN<sup>+</sup>21]; the complete  $\text{rwPHP(PLS)}$ -hardness result can be found in Section 5.2.

There is a clear intuition behind the reduction: suppose  $\text{rwPHP}$  were false, i.e., there are functions  $f : [N] \rightarrow [2N]$  and  $g : [2N] \rightarrow [N]$  such that  $f \circ g : [2N] \rightarrow [2N]$  is the identity function, then every unsatisfiable CNF  $F$  would have a resolution refutation of size  $\text{poly}(N, n)$ . Of course, the ground truth is that such functions  $f$  and  $g$  should not exist, but a weak proof system (such as resolution itself) might not be aware of this. Suppose the weak system “thinks” that such a pair of functions  $(f, g)$  might exist, and it can construct a short resolution refutation of  $F$  from  $(f, g)$ , then the weak system should also “think” that  $F$  might have a short resolution refutation. In summary, if it is hard to refute the existence of  $(f, g)$  (which means proving  $\text{rwPHP}$ ), then it is also hard to prove that  $F$  does not have a short resolution refutation.

Now, our task becomes the following. We live in a strange world where there is a surjection from  $[N]$  to  $[2N]$ ; given an arbitrary unsatisfiable CNF  $F$ , we want to construct a  $\text{poly}(N, n)$ -size resolution refutation of  $F$ . Consider the size- $2^{O(n)}$  brute-force resolution refutation for every unsatisfiable CNF, which is represented by the following proof tree.

- The root (level 0) of the tree contains the empty clause  $\perp$ .
- For each level  $1 \leq i \leq n$ , each clause  $C$  at level  $i - 1$  is resolved from the two clauses  $C \vee x_i$  and

$C \vee \bar{x}_i$ , both of which sits in level  $i$ . The clauses  $C \vee x_i$  and  $C \vee \bar{x}_i$  are the two *children* of  $C$ . Note that each clause at level  $i$  has a width of exactly  $i$ .

- Finally, every clause  $C$  at level  $n$  corresponds to an assignment  $x_C \in \{0, 1\}^n$  which is the only assignment falsifying  $C$ . Since  $F$  is unsatisfiable, there is an axiom of  $F$  that  $x_C$  falsifies. Clearly,  $C$  is a *weakening* of this axiom.

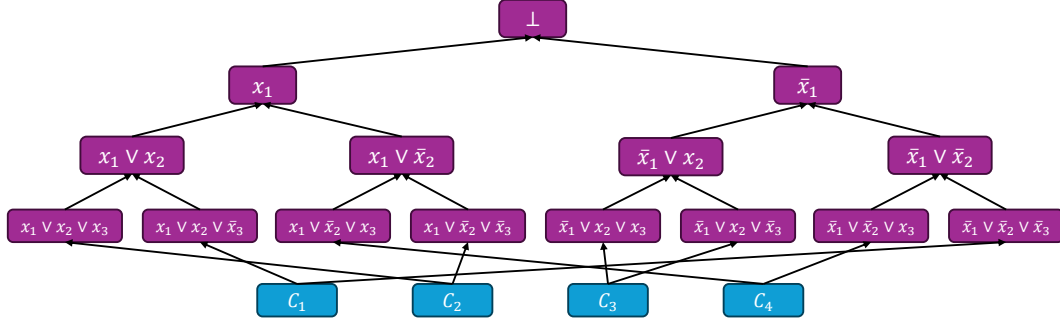


Figure 3: The brute-force resolution proof for a CNF  $F = C_1 \wedge C_2 \wedge C_3 \wedge C_4$  when  $n = 3$ .

We now construct a shorter resolution refutation using the surjection from  $[N]$  to  $[2N]$ . We guarantee that in our short refutation, each level never contains more than  $N$  clauses; this implies that our resolution refutation is of size  $O(N \cdot n)$ . Consider level  $i$  where  $1 \leq i \leq n$ . If level  $i - 1$  contains at most  $N$  clauses, then level  $i$  contains at most  $2N$  clauses: for each clause  $C_j$  in level  $i - 1$ , there are two clauses  $C'_{2j} := C \vee x_i$  and  $C'_{2j+1} := C \vee \bar{x}_i$  in level  $i$ . However, since there is a *surjection* from  $[N]$  to  $[2N]$ , it is possible to pick  $N$  clauses among these  $2N$  ones such that each of the  $2N$  clauses appears in these  $N$  ones! (The  $j$ -th picked clause ( $j \in [N]$ ) is  $C'_{f(j)}$ ; the clause  $C'_j$  ( $j \in [2N]$ ) appears as the  $g(j)$ -th picked clause.) Now that level  $i$  also contains at most  $N$  clauses, we can proceed to the next level and so on.

We stress again that the ground truth is, of course, that there do not exist functions  $f : [N] \rightarrow [2N]$  and  $g : [2N] \rightarrow [N]$  such that  $f \circ g : [2N] \rightarrow [2N]$  is the identity function. However, the point is that given any step in the above resolution refutation that is an invalid derivation, we can pinpoint a “witness” number  $x \in [2N]$  such that  $f(g(x)) \neq x$ .

The above describes the intuition behind the decision tree reduction from rwPHP to the refuter problems of resolution size lower bounds presented in [dRGN<sup>+</sup>21]. Our reduction from rwPHP(PLS) to the refutation problems proceeds in the same way, except that now  $g$  is only a function computable in PLS. Compared with [AM20, dRGN<sup>+</sup>21], our proof only has one more component: showing that these PLS instances can also be embedded into the above resolution refutation. We refer the reader to the formal proof in Section 5.2 for details.

### 3 Preliminaries

The first three subsections present standard preliminaries and can be skipped if the reader is familiar. However, the last two subsections introduce new concepts and it is highly recommended to read through (i.e., not skip) them. In particular, Section 3.4 introduces the refuter problems for resolution lower bounds as TFNP<sup>dt</sup> problems, and Section 3.5 defines and discusses the subclass rwPHP(PLS).

We use 0-indexing:  $[n] = \{0, 1, \dots, n - 1\}$ . For functions  $f : \mathcal{A} \rightarrow \mathcal{B}$  and  $g : \mathcal{B} \rightarrow \mathcal{C}$ , their *composition*  $g \circ f$  is defined as

$$\forall x \in \mathcal{A}, (g \circ f)(x) = g(f(x)).$$

### 3.1 Pigeonhole Principle

Let  $m > n$ , the *pigeonhole principle* (PHP) states that there is no way to send  $m$  pigeons into  $n$  holes such that different pigeons are sent to different holes. This is expressed as the following unsatisfiable CNF  $\text{PHP}_{m \rightarrow n}$ . (In the definition below, think of  $x_{ij} = 1$  if pigeon  $i$  goes to hole  $j$ .)

**Definition 3.1** ( $\text{PHP}_{m \rightarrow n}$ ).  $\text{PHP}_{m \rightarrow n}$  is the conjunction of the following set of clauses:

- $\bigvee_{j \in [n]} x_{ij}$  for every pigeon  $i \in [m]$ ;
- $\bar{x}_{ij} \vee \bar{x}_{i'j}$  for every two different pigeons  $0 \leq i < i' \leq m - 1$  and every hole  $j \in [n]$ .

The seminal work of Haken [Hak85] proved that any resolution proof of  $\text{PHP}_{(n+1) \rightarrow n}$  requires  $2^{\Omega(n)}$  size. The proof of this classical theorem has been simplified by several follow-up works [CP90, BP96, BW01].

### 3.2 Decision Tree TFNP

Let  $\mathcal{O} = \{O_N\}_N$  be a family of solution spaces. A *search* problem  $\mathcal{P}$  is a family of sets  $\{P_N\}_N$ , where each  $P_N$  is a subset of  $\{0, 1\}^N \times O_N$ . Let  $x \in \{0, 1\}^N$  be an *input* to  $\mathcal{P}$ , we say that  $o \in O_N$  is a *solution* of  $x$  if  $(x, o) \in P_N$ . We say  $\mathcal{P}$  is *total* if every  $x \in \{0, 1\}^*$  has at least one solution. We sometimes abuse the notation by calling an individual relation  $P_N$  a search problem, and implicitly assume that there is a sequence  $\{P_N\}_N$ .

We study total search problems in the *decision tree* model. In this model, we think of the input  $x \in \{0, 1\}^N$  as very long and can only be accessed by querying individual bits. An algorithm (i.e., decision tree) is *efficient* if it only makes  $\text{polylog}(N)$  many queries. We will typically consider search problems where  $|O_N| \leq 2^{\text{polylog}(N)}$ , so efficient algorithms will be able to handle solutions  $o \in O_N$  in their entirety. A search problem  $\mathcal{P}$  is in  $\text{FNP}^{\text{dt}}$  if given (oracle access to) an input  $x \in \{0, 1\}^N$  and a solution  $o \in O_N$ , there is an efficient decision tree  $T_o$  for deciding whether  $(x, o) \in P_N$ . The class  $\text{TFNP}^{\text{dt}}$  consists of all *total* search problems in  $\text{FNP}^{\text{dt}}$ .

For example, an important  $\text{TFNP}^{\text{dt}}$  problem in this paper is the problem ITER, defined as follows.

#### Problem ITER

Input: A function  $S : [N] \rightarrow [N]$ .

Output: A number  $x \in [N]$  is a valid solution if one of the following holds:

- $x = 0$  and  $S(0) = 0$ ;
- $S(x) < x$ ; or
- $S(x) > x$  and  $S(S(x)) = S(x)$ .

It is easy to check that ITER is in  $\text{FNP}^{\text{dt}}$ : Given an output  $x$  and oracle access to the function  $S : [N] \rightarrow [N]$ , one can verify whether  $x$  is a valid solution by querying at most 2 entries of  $S$ ; namely  $S(x)$  and  $S(S(x))$ . Since each entry can be represented by at most  $\log N$  bits, the query complexity of verifying solutions for ITER is  $\text{polylog}(N)$ . On the other hand, the totality of ITER expresses the following fact: every DAG has a sink. It turns out that we will also frequently use a reversed version of ITER for simplicity, whose equivalence to ITER is easy to see: Given a function  $S : [N] \rightarrow [N]$  such that  $S(N - 1) < N - 1$ , find some  $x \in [N]$  such that 1) either  $S(x) > x$  or 2)  $S(x) < x$  and  $S(S(x)) = S(x)$ .

**Definition 3.2** (Decision tree reductions). Let  $\mathcal{P}, \mathcal{Q}$  be two  $\text{TFNP}^{\text{dt}}$  problems, and  $d(N)$  be a parameter (typically  $\text{polylog}(N)$ ). A *depth- $d$*  decision tree reduction from  $\mathcal{P}$  to  $\mathcal{Q}$  consists of two functions  $(f, g)$ , where each output bit of  $f, g$  can be computed from the input  $x$  by a depth- $d$  decision tree:

- $f : \{0, 1\}^N \rightarrow \{0, 1\}^{M(N)}$  maps an input  $x$  of  $\mathcal{P}$  to an input  $f(x)$  of  $\mathcal{Q}$ .

- $g$  maps any valid solution of  $f(x)$  (as an instance of  $\mathcal{Q}$ ) into a valid solution of  $x$  (as an instance of  $\mathcal{P}$ ).

We say the reduction is *uniform* if both  $f$  and  $g$  can be computed by uniform Turing machines with query access to  $x$ . We allow  $M(N)$  to be super-polynomial in  $N$ , but we require  $M(N) \leq \exp(d(N))$ .

Usually, for two  $\text{TFNP}^{\text{dt}}$  problems  $\mathcal{P}, \mathcal{Q}$  we say  $\mathcal{P}$  can be (*many-one*) reduced to  $\mathcal{Q}$  if there is a  $\text{polylog}(N)$ -depth decision tree reduction from  $\mathcal{P}$  to  $\mathcal{Q}$ .

The class  $\text{PLS}^{18}$  is the class of problems in  $\text{TFNP}^{\text{dt}}$  that has a depth- $\text{polylog}(N)$  reduction to ITER. (Note that PLS was originally defined differently [JPY88]; the PLS-completeness of ITER was shown in [Mor01].)

The inputs of most  $\text{TFNP}^{\text{dt}}$  problems introduced in this paper will be partitioned into blocks; for example, the input of ITER consists of  $N$  blocks where each block consists of  $\log N$  bits describing an integer in  $[N]$ . It will be more convenient to work with the *block-depth* of decision trees, which is the number of different *blocks* that a decision tree queries. For example, solutions of ITER can be verified in block-depth 2. The problems in this paper will have block size  $\text{polylog}(N)$ , hence  $\text{polylog}(N)$  block-depth is equivalent to  $\text{polylog}(N)$  (bit-)depth. However, we will upper bound the complexity of our decision trees by block-depth for convenience. Although the distinction of depth and block-depth does not make an essential difference in this paper, many interesting lifting theorems and non-automatability results are recently proved using the notion of block-depth (or block-width) [AM20, GKMP20, dRGN<sup>+</sup>21]. It might be beneficial to have bounds on block-depth, which is usually sharper as the decision trees we construct tend to query many bits in the same block.

We assume all the  $\text{TFNP}^{\text{dt}}$  problems discussed in this paper are *paddable*, i.e., for any  $N < M$ , solving an instance of size  $N$  could always be efficiently reduced to solving an instance of size  $M$  of the same problem. Most of the common  $\text{TFNP}^{\text{dt}}$  problems can be easily formulated in a paddable way.<sup>19</sup>

### 3.2.1 Connection to Proof Complexity

There is a generic connection between  $\text{TFNP}^{\text{dt}}$  and propositional proof complexity via the *false clause search* problem (see, e.g. [dRGR22, BFI23]).

**Definition 3.3.** For an unsatisfiable CNF  $F := C_1 \wedge \dots \wedge C_m$ ,  $\text{Search}(F)$  is the search problem in which an assignment  $x$  to  $F$  is given via query access, and a solution is a clause  $C_i$  of  $F$  falsified by  $x$ .

Define  $\text{Search}(\mathcal{F})$  for a family of formula  $\mathcal{F} = \{F_n\}_{n \in \mathbb{N}}$  as  $\{\text{Search}(F_n)\}_{n \in \mathbb{N}}$  accordingly.

When the width of  $F$  is  $\text{polylog}(n)$ , where  $n$  is the number of variables in  $F$ ,  $\text{Search}(F)$  is a  $\text{TFNP}^{\text{dt}}$  problem. In the other direction, for any  $\text{TFNP}^{\text{dt}}$  problem  $R_n \in \{0, 1\}^n \times O_n$ , it can be equivalently written as  $\text{Search}(F_n)$  for some CNF  $F_n$  of  $\text{polylog}(n)$  width. More specifically, let  $\{T_o\}_{o \in O_n}$  be the set of efficient decision trees for verifying solutions,  $\neg T_o(x)$  can be written as a low-width CNF stating that any accepting path in  $T_o$  is falsified by  $x$ . We then take

$$F_n = \bigwedge_{o \in O_n} \neg T_o(x), \tag{2}$$

and it is easy to see the equivalence between  $\text{Search}(F_n)$  and  $R_n$  by definition.

Informally, we say a proof system  $P$  is characterized by a syntactical  $\text{TFNP}^{\text{dt}}$  subclass  $\mathcal{C}$  if for any family of formula  $\mathcal{F} = \{F_n\}$ ,  $P$  has a *small* proof of  $\mathcal{F}$  if and only if  $\text{Search}(\mathcal{F}) \in \mathcal{C}$ . Buss, Fleming,

<sup>18</sup>In this paper, most of the times when we mention a syntactic subclass of  $\text{TFNP}$  (such as  $\text{PLS}$ ) we mean the decision tree version of it (i.e.,  $\text{PLS}^{\text{dt}}$ ), and it should be easy to figure out whether we mean the decision tree version or the Turing machine version of this subclass from the context. Therefore, for convenience, we drop the superscript  $\text{dt}$  when we express syntactic subclasses of  $\text{TFNP}^{\text{dt}}$ . We still preserve the superscript  $\text{dt}$  in “ $\text{TFNP}^{\text{dt}}$ ” when we want to emphasize that the underlying model is decision tree  $\text{TFNP}$ .

<sup>19</sup>There is a similar notion called *instance extension*, which is defined in [BM04].

and Impagliazzo [BFI23] showed that any well-behaved<sup>20</sup> proof system  $P$  is characterized by a TFNP<sup>dt</sup> subclass  $C$ , and vice versa. In particular, resolution is characterized by PLS.

**Theorem 3.4** (Folklore). *Let  $\mathcal{F} = \{F_n\}$  be a family of unsatisfiable formula,  $\text{Search}(\mathcal{F}) \in \text{PLS}$  if and only if  $\mathcal{F}$  have a polylog( $n$ )-width resolution refutation.*

### 3.3 Bounded Arithmetic

We introduce the theories  $\mathsf{T}_2^1$ ,  $\mathsf{T}_2^1 + \text{dwPHP}(\text{PV})$ , as well as their relativized versions. A more comprehensive introduction of bounded arithmetic (including the theories  $\mathsf{S}_2^i$  and  $\mathsf{T}_2^i$ ) can be found in [Kra95].

The language of bounded arithmetic consists of the following symbols

$$L_{\text{BA}} := \{0, 1, +, \cdot, <, =, \lfloor \cdot / 2 \rfloor, |\cdot|, \#\}.$$

Here, the intended meaning of  $|a|$  is the *bit-length* of the binary number  $a$ , i.e.,

$$|a| := \begin{cases} \lceil \log_2(a+1) \rceil & \text{if } a > 0; \\ 0 & \text{if } a = 0. \end{cases}$$

The intended meaning of  $\#$  (“smash”) is

$$x \# y := 2^{|x| \cdot |y|};$$

roughly speaking, this symbol is used to create objects whose size is *polynomial*, instead of only *linear*, in the length of its inputs. These symbols are governed by a list of 32 axioms called BASIC, each of which asserts some basic fact about the intended meanings of these symbols. For instance:

$$a \leq b \rightarrow a \leq b + 1. \quad (\text{axiom 1 in BASIC})$$

The complete list of BASIC axioms can be found in [Kra95, Definition 5.2.1].

A *bounded quantifier* is a quantifier of the form

$$\forall y < t(\vec{x}) \quad \text{or} \quad \exists y < t(\vec{x})$$

for some term  $t$ . Formally, they are defined as abbreviations:

$$\begin{aligned} \forall y < t(\vec{x}) \varphi(\vec{x}, y) &:= \forall y (y < t(\vec{x}) \rightarrow \varphi(\vec{x}, y)); \\ \exists y < t(\vec{x}) \varphi(\vec{x}, y) &:= \exists y (y < t(\vec{x}) \wedge \varphi(\vec{x}, y)). \end{aligned}$$

A *sharply bounded quantifier* is a quantifier of the form

$$\forall y < |t(\vec{x})| \quad \text{or} \quad \exists y < |t(\vec{x})|.$$

That is, the domain of possible values of  $y$  is bounded by the length of a term. Intuitively, sharply bounded quantifiers are “feasible” because, thinking of  $t(\vec{x})$  as the description of a polynomial-size object, there are only polynomially many possibilities of  $y$  and they can be enumerated in polynomial time.

A formula is *sharply bounded* if all quantifiers in it are sharply bounded quantifiers. A  $\Sigma_1^b$ -formula is a formula constructed from sharply bounded formulas using  $\wedge$ ,  $\vee$ , sharply bounded quantifiers, and *existential bounded quantifiers* (“ $\exists y < t(\vec{x})$ ”). It can be shown that the languages defined by  $\Sigma_1^b$ -formulas are exactly those computed in NP.

The power of theories in bounded arithmetic comes from their *induction axioms*. Let  $\Phi$  be a class of formulas, then  $\Phi\text{-IND}$  is the following axiom schema

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

<sup>20</sup>Here, a proof system is well-behaved if it is closed under decision tree reduction, and it can prove its own soundness.

for every  $\phi \in \Phi$ . The definition of  $\mathsf{T}_2^1$  is:

$$\mathsf{T}_2^1 := \text{BASIC} + \Sigma_1^b\text{-IND.}$$

That is, when reasoning in  $\mathsf{T}_2^1$ , it is allowed to use induction axioms over  $\Sigma_1^b$  formulas (i.e., NP languages).

It is equivalent, and sometimes more convenient to replace  $\Sigma_1^b\text{-IND}$  with  $\Sigma_1^b\text{-MIN}$ , the *minimization principle* over  $\Sigma_1^b$  formulas. For a set of formulas  $\Phi$ , the axiom schema  $\Phi\text{-MIN}$  consists of

$$\phi(a) \rightarrow \exists x \leq a \forall y < x (\phi(x) \wedge \neg\phi(y))$$

for every  $\phi \in \Phi$ . Equivalently, when reasoning in  $\mathsf{T}_2^1$ , it is allowed to use the fact that there exists a *smallest*  $x$  such that  $C(x) = 1$ , whenever  $C$  is a polynomial-size *nondeterministic circuit* and we know some  $y$  such that  $C(y) = 1$ .

The theory  $\text{PV}$  is an equational theory defined by Cook [Coo75] to capture polynomial-time reasoning. It contains a function symbol for every polynomial-time algorithm, introduced inductively using Cobham's recursion-theoretic characterization of polynomial time [Cob64]. More detailed treatments about  $\text{PV}$  can be found in [Kra95, CN10, CLO24]. In the literature, it is common to also use  $\text{PV}$  to denote the set of function symbols in  $\text{PV}$  (which corresponds to functions computable in polynomial time).

The *dual weak pigeonhole principle* over  $\text{PV}$  functions, denoted as  $\text{dwPHP}(\text{PV})$ , is the following axiom schema

$$\forall a > 1 \exists v < a^2 \forall u < a f(u) \neq v$$

for every  $\text{PV}$ -function  $f$  with parameters<sup>21</sup>. Roughly speaking, this means that if we have a polynomial-size circuit  $f : \{0, 1\}^n \rightarrow \{0, 1\}^{2n}$  (think of  $a = 2^n$  above), then there exists some  $v \in \{0, 1\}^{2n}$  that is not in the range of  $C$ . We note that the choice of  $a^2$  above is somewhat arbitrary, as  $\text{dwPHP}(\text{PV})$  with various parameters are equivalent over  $\mathsf{S}_2^1 \subseteq \mathsf{T}_2^1$  [PWW88, Jeř04].

To summarize, when reasoning in the theory  $\mathsf{T}_2^1 + \text{dwPHP}(\text{PV})$ , one is allowed to use the following two axiom schemas:

( $\Sigma_1^b\text{-MIN}$ ) For a polynomial-size nondeterministic circuit  $C$  and some  $y$  such that  $C(y) = 1$ , there exists a *smallest*  $x$  such that  $C(x) = 1$ .

( $\text{dwPHP}(\text{PV})$ ) For a polynomial-size circuit  $C : \{0, 1\}^n \rightarrow \{0, 1\}^{2n}$ , there exists a string  $y \in \{0, 1\}^{2n}$  that is not in the range of  $C$ .

Finally, the *relativized* theories  $\mathsf{T}_2^1(\alpha)$  and  $\mathsf{T}_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha))$  are simply their unrelativized counterparts with a new unary relation symbol  $\alpha$  added into the language  $L_{\text{BA}}$ . (One can think of  $\alpha$  as an oracle that encodes an exponentially-long input; for example,  $\alpha(i)$  might encode the  $i$ -th bit of an exponentially-long resolution proof according to some canonical encoding.) The class of  $\Sigma_1^b(\alpha)$  formulas and axioms  $\Sigma_1^b(\alpha)\text{-MIN}$  and  $\text{dwPHP}(\text{PV}(\alpha))$  are relativized in a straightforward way. There are no other axioms involving  $\alpha$  except for the induction axioms and dual weak pigeonhole principles. To summarize:

- When reasoning in  $\mathsf{T}_2^1(\alpha)$ , it is allowed to use  $\Sigma_1^b(\alpha)\text{-MIN}$ , i.e., for any polynomial-size nondeterministic *oracle* circuit  $C^\alpha$  and input  $y$  such that  $C^\alpha(y) = 1$ , there exists a *smallest* input  $x$  such that  $C^\alpha(x) = 1$ .
- When reasoning in  $\mathsf{T}_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha))$ , it is additionally allowed to use the fact that for any polynomial-size *oracle* circuit  $C^\alpha : \{0, 1\}^n \rightarrow \{0, 1\}^{2n}$ , there exists some  $y \in \{0, 1\}^{2n}$  that is not in the range of  $C^\alpha$ .

---

<sup>21</sup>This is the standard terminology in bounded arithmetic that means  $f$  might depend on some other parameter not shown above. The parameter can be thought of as non-uniformity; cf. [Footnote 11](#).

### 3.4 Refuter Problems for Resolution Lower Bounds

We provide formal definitions of the refuter problems in the decision tree model. We begin by defining *resolution refutations*; the definition is adapted from [dRGN<sup>+</sup>21, Section 3.1].

**Definition 3.5.** Let  $F$  be an unsatisfiable CNF with  $n$  variables and  $m$  clauses; the clauses in  $F$  will be called *axioms* and will be denoted as  $C_{-m}, \dots, C_{-1}$  for convenience. A *resolution refutation* of  $F$  is a sequence of nodes  $C_0, C_1, \dots, C_{L-1}$ , where each node  $C_i$  contains the following information.

- A set of literals among  $\{x_1, x_2, \dots, x_n, \bar{x}_1, \bar{x}_2, \dots, \bar{x}_n\}$ . Abusing notation, we also denote the clause consisting of the disjunction of these literals by  $C_i$ .
- A *tag* which is one of the following: “resolution” or “weakening”.
- Two integers  $-m \leq j, k < i$  and a variable  $a \in \{1, 2, \dots, n\}$  if the tag is “resolution”. This means that  $C_i$  is obtained from the clauses  $C_j$  and  $C_k$  by *resolving* the variable  $x_a$ .
- One integer  $-m \leq j < i$  if the tag is “weakening”. This means that  $C_i$  is a *weakening* of  $C_j$ .

The resolution refutation is valid if the following is true for every  $1 \leq i \leq L$ :

- If  $C_i$  is marked “resolution”, then there are clauses  $D$  and  $E$  such that  $C_j = x_a \vee D$ ,  $C_k = \bar{x}_a \vee E$ , and  $C_i = D \vee E$ .
- If  $C_i$  is marked “weakening”, then there is a clause  $D$  such that  $C_i = C_j \vee D$ .
- Finally,  $C_{L-1} = \perp$  (i.e., contains no literals).

The *length* or *size* of the refutation is  $L$ , and the *width* of the refutation is the maximum integer  $w$  such that every clause  $C_i$  ( $-m \leq i < L$ ) in the refutation contains at most  $w$  literals.

Resolution is *complete* and *sound*: a CNF  $F$  has a resolution refutation (of whatever length) if and only if it is unsatisfiable.

Each node in the resolution refutation would be a *block*; therefore, when we say a decision tree over a resolution refutation has block-depth  $d$ , we mean that it only queries (potentially all information in)  $d$  nodes of the refutation.

Next, we define the refuter problems.

**Definition 3.6.** Let  $\mathcal{F} = \{F_n\}_{n \in \mathbb{N}}$  be a family of unsatisfiable CNFs where every  $F_n$  requires resolution of width greater than  $w_n$  and size greater than  $s_n$ .

- An input to the problem  $\text{REFUTER}(w(F_n \vdash_{\text{Res}} \perp) \leq w_n)$  is a purported resolution refutation of  $F_n$  with width at most  $w_n$ . (It is easy to syntactically guarantee that the width of the input refutation is at most  $w_n$  by allocating only  $w_n$  literals for each node.)
- An input to the problem  $\text{REFUTER}(s(F_n \vdash_{\text{Res}} \perp) \leq s_n)$  is a purported resolution refutation of  $F$  with at most  $s_n$  clauses.

The outputs of these problems consist of only one index  $i$ , which means the node  $C_i$  does not satisfy the validity conditions defined in [Definition 3.5](#). We will call such nodes *invalid derivations* or *illegal derivations*.

Note that each node can be described in  $\text{poly}(w, \log n, \log L)$  bits where  $w$  is the width of the resolution refutation. Hence, in the typical parameter regime, we will consider resolution refutations whose length is exponential in its width ( $L = 2^{w^{\Omega(1)}}$ ), so that the access to each block is “efficient”, i.e., only needs to query polylogarithmic many bits. In particular, the typical parameter regime for size lower bounds is exponential

( $L = 2^{n^{\Omega(1)}}$ ); polynomial width lower bound ( $w = n^{\Omega(1)}$ ) is considered in [Section 4](#) and [Section 6](#), while polylogarithmic width lower bound is considered in [Section 7.2](#). We also assume  $L = 2^{n^{O(1)}}$ , so that the proof is not extremely redundant.

Since there is a decision tree of block-depth at most 3 verifying whether a given node in the resolution refutation is an invalid derivation, the refuter problems defined above are in  $\text{FNP}^{\text{dt}}$ . Moreover, if the resolution lower bounds ( $w_0$  or  $s_0$ ) are indeed true, then the refuter problems defined above are total. Hence,  $\text{REFUTER}(\cdot)$  is a natural family of problems in  $\text{TFNP}^{\text{dt}}$ .

### 3.5 $\mathcal{P}$ -Retraction Weak Pigeonhole Principle

Recall that  $\text{rwPHP}$ , the *retraction weak pigeonhole principle*, is the following principle:

**Fact 3.7.** *Let  $g : [2M] \rightarrow [M]$  and  $f : [M] \rightarrow [2M]$  be two functions. Then there must exist some  $y \in [2M]$  such that  $f(g(y)) \neq y$ .*

Roughly speaking, for any TFNP class  $\mathcal{P}$ ,  $\text{rwPHP}(\mathcal{P})$  is the retraction weak pigeonhole principle where the retraction ( $g : [2M] \rightarrow [M]$ ) is a (multi-valued) function computable in  $\mathcal{P}$ . For example:

#### Problem $\text{rwPHP}(\text{PLS})$

Input: Let  $M \leq N/2$ . The input consists of the following functions:

- $f : [M] \rightarrow [N]$  is a purported “surjection”;
- for each  $y \in [N]$ ,  $I_y := (L, S_y)$  is an instance of  $\text{ITER}$ , where  $S_y : [L] \rightarrow [L]$ ; and
- $g_y : [L] \rightarrow [M]$  maps solutions of  $I_y$  to integers in  $[M]$ .

Output: A number  $y \in [N]$  and  $ans \in [L]$  such that  $ans$  is a solution of the  $\text{ITER}$  instance  $g_y$  and  $f(g_y(ans)) \neq y$ .

In general, for a  $\text{TFNP}^{\text{dt}}$  problem  $\mathcal{P}$ , we define  $\text{rwPHP}(\mathcal{P})$  by replacing each  $I_y$  in the above definition with an instance of  $\mathcal{P}$ . For convenience, we sometimes use the same notation  $\text{rwPHP}(\mathcal{P})$  to indicate the TFNP subclass that contains any TFNP problem reducible to the problem  $\text{rwPHP}(\mathcal{P})$ . It will be clear from the context whether  $\text{rwPHP}(\mathcal{P})$  refers to the problem or the class in the remainder of the paper.

**Fact 3.8.**  $\text{rwPHP}(\mathcal{P}) \in \text{TFNP}^{\text{dt}}$ .

*Proof.* To verify a solution  $(y, ans)$ , check that  $ans$  is a valid solution for  $I_y$  and that  $f(g_y(ans)) \neq y$ .

The totality (i.e., existence of solutions) can be argued as follows. For each  $y \in [N]$ , let  $ans(y)$  be a solution (say the lexicographically first one) of  $I_y$ ; since  $\mathcal{P}$  is a total problem,  $ans(y)$  exists. Let  $g'(y) := g_y(ans(y))$ . By the retraction weak pigeonhole principle, there exists some  $y \in [N]$  such that  $f(g'(y)) \neq y$ . It follows that  $(y, ans(y))$  is a valid solution for  $\text{rwPHP}(\mathcal{P})$ .  $\square$

**Fact 3.9.** *There is a depth-1 decision tree reduction from  $\mathcal{P}$  to  $\text{rwPHP}(\mathcal{P})$  and a depth-1 decision tree reduction from  $\text{rwPHP}$  to  $\text{rwPHP}(\mathcal{P})$ .*

*Proof.* To reduce  $\mathcal{P}$  to  $\text{rwPHP}(\mathcal{P})$ : let  $I$  be a  $\mathcal{P}$  instance. Define  $f(x) = 1$  as a trivial function; for each  $y \in [N]$ , define the instance  $I_y := I$ ; for every possible answer  $ans$  of  $I_y$ , let  $g_y(ans) = 1$ . Clearly, for any answer  $(y, ans)$  of the  $\text{rwPHP}(\mathcal{P})$  instance,  $ans$  itself would be a valid answer of the  $\mathcal{P}$  instance  $I$ .

To reduce  $\text{rwPHP}$  to  $\text{rwPHP}(\mathcal{P})$ : let  $f : [M] \rightarrow [N]$  and  $g : [N] \rightarrow [M]$  be an  $\text{rwPHP}$  instance. Fix any (say trivial)  $\mathcal{P}$  instance  $I$ . For each  $y \in [N]$ , define the instance  $I_y := I$ ; for every possible answer  $ans$  of  $I_y$ , let  $g_y(ans) = g(y)$ . For any answer  $(y, ans)$  of the  $\text{rwPHP}(\mathcal{P})$  instance, since  $f(g_y(ans)) \neq y$ , it follows that  $f(g(y)) \neq y$ , and hence  $y$  is a valid answer of the  $\text{rwPHP}$  instance  $(f, g)$ .  $\square$

**Fact 3.10.**  $\text{rwPHP}(\text{PLS})$  is not in  $\text{PLS}$  in the black-box setting.

*Proof.* Classical techniques (such as Prover-Delayer games) show that the totality of rwPHP requires resolution width  $\Omega(M)$  to prove (see also [PT19, dRGN<sup>+</sup>21]), which means that any decision tree of depth  $o(M)$  cannot reduce rwPHP to PLS. It follows from Fact 3.9 that there is a black-box separation between rwPHP(PLS) and PLS itself.  $\square$

**Fact 3.11.** *Let  $\mathcal{P}$  and  $\mathcal{Q}$  be TFNP<sup>dt</sup> problems. If there is a depth- $d$  decision tree reduction from  $\mathcal{P}$  to  $\mathcal{Q}$ , then there is a depth- $d$  decision tree reduction from rwPHP( $\mathcal{P}$ ) to rwPHP( $\mathcal{Q}$ ).*

*Proof.* Let  $f : [M] \rightarrow [N]$ ,  $\{I_y\}_{y \in [N]}$ , and  $\{g_y\}_{y \in [N]}$  be an instance of rwPHP( $\mathcal{P}$ ). Define an instance  $(f', \{I'_y\}, \{g'_y\})$  of rwPHP( $\mathcal{Q}$ ) as follows. The function  $f' := f$  stays the same; each  $I'_y$  is obtained by running the reduction from  $\mathcal{P}$  to  $\mathcal{Q}$  on  $I_y$ ; for each possible answer  $ans'$  of each  $I'_y$ , we run the reduction to obtain an answer  $ans$  of  $I_y$ , and return  $g'_y(ans') := g_y(ans)$ .

Finally, let  $(y', ans')$  be a valid answer of the instance  $(f', \{I'_y\}, \{g'_y\})$ . Let  $y := y'$  and run the reduction to obtain an answer  $ans$  of  $I_y$  from the answer  $ans'$  of  $I'_y$ , then  $(y, ans)$  is a valid answer of  $(f, \{I_y\}, \{g_y\})$ .  $\square$

It follows from Fact 3.11 that, for example, the class rwPHP(PLS) can be defined from *any* complete problem for PLS.

**Amplification for rwPHP.** In this paper, when we talk about rwPHP, we always think about a purported “surjection”  $f : [M] \rightarrow [N]$  where  $N = 2M$ . This is without loss of generality since the complexity of rwPHP does not depend significantly on the relationship between  $M$  and  $N$  (unless they are too close to each other). This is also true for rwPHP( $\mathcal{P}$ ) provided that  $\mathcal{P}$  is closed under *Turing reductions*. We note that many interesting subclasses of TFNP (such as PLS, PPA, PPAD, and PPADS) are indeed closed under Turing reductions [BJ12, Section 6], with the notable exception of PPP in the black-box model [FGPR24].

**Theorem 3.12** (Informal). *Suppose  $\mathcal{P}$  is closed under Turing reductions. Let  $\text{rPHP}_{M \rightarrow N}(\mathcal{P})$  denote the problem rwPHP( $\mathcal{P}$ ) where the given “purported surjection” is from  $[M]$  to  $[N]$ .<sup>22</sup> Then there is an efficient decision tree reduction from  $\text{rPHP}_{M \rightarrow (M+M/\text{polylog}(M))}(\mathcal{P})$  to  $\text{rPHP}_{M \rightarrow M^{100}}(\mathcal{P})$ .*

We remark that amplification of weak pigeonhole principles is a well-known fact in bounded arithmetic [PWW88, Tha02, Kra04, Jeř04, Jeř07b, CLO24] and total search problems [Kor21, Kor22]. Since the proof follows from standard arguments in the literature, we postpone it to Appendix A.

### 3.5.1 Witnessing for $\text{T}_2^1 + \text{dwPHP}(\text{PV})$

As mentioned in the introduction, every TFNP problems whose totality can be proved in  $\text{T}_2^1 + \text{dwPHP}(\text{PV})$  reduces to rwPHP(PLS). This is an easy corollary of [BK94] and [AT14, Lemma 2.1]; for completeness we present a proof here.

**Theorem 3.13.** *Suppose that  $\phi(x) := \exists y < t_\phi(x) \psi(x, y)$  is a  $\Sigma_1^b(\alpha)$ -sentence and*

$$\text{T}_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha)) \vdash \forall x \phi(x).$$

*Then the TFNP problem corresponding to  $\phi$  is in rwPHP(PLS).*

*Proof.* It is shown in [AT14, Lemma 2.1] that there is a term  $t = t(x)$  and a function symbol  $f \in \text{PV}(\alpha)$  such that

$$\text{T}_2^1(\alpha) \vdash \forall x (t > 2 \wedge \forall v < t^2 \exists u < t (f_x(u) = v \vee \phi(x))). \quad (3)$$

<sup>22</sup>The notation  $\text{rPHP}_{M \rightarrow N}$  means “retraction pigeonhole principle from  $M$  pigeons to  $N$  holes”. As “weak” conventionally refers to the case where  $M = 2N$ , the retraction pigeonhole principle with  $M$  and  $N$  specified are called  $\text{rPHP}_{M \rightarrow N}$  instead of  $\text{rPHP}_{M \rightarrow N}$ .

(In fact, this follows from standard manipulations underlying Wilkie’s witnessing theorem for  $S_2^1 + \text{dwPHP}(\text{PV})$ ; see e.g., [Jef04, Proposition 14].)

To parse (3), note that  $f_x : [t] \rightarrow [t^2]$  defines a “stretching” function and hence cannot be surjective. The non-existence of  $u$  such that  $f_x(u) = v$  exactly means that  $v$  is not in the range of  $f_x$ ; hence given (3),  $\forall x \phi(x)$  follows from the dual weak pigeonhole principle.

The following problem is in PLS by the witnessing theorem for  $T_2^1$  [BK94]. On input  $(x, v)$ , find either some  $u < t(x)$  such that  $f_x(u) = v$  or some  $y < t_\phi(x)$  such that  $\psi(x, y)$  holds (we call such  $y$  a “certificate” for  $x$ ). Now it is at least easy to see that there is a *randomized* reduction from the TFNP problem corresponding to  $\phi(x)$  to PLS: Let  $v \leftarrow [t^2]$  be random (which is a non-output of  $f_x$  w.h.p.), then the above PLS procedure finds some  $y$  that is a certificate for  $x$ .

Working slightly harder, we can see that the above is actually a reduction to  $\text{rwPHP}(\text{PLS})$ :

- The “purported surjection” is the function  $f_x : [t] \rightarrow [t^2]$ .
- For each  $v \in [t^2]$ , there is a PLS instance  $I_v$  which captures the problem of given  $(x, v)$  outputting either  $u \in f_x^{-1}(v)$  or a certificate  $y$  for  $x$ .

Hence, given any  $v \in [t^2]$  and solution  $ans$  of  $I_v$  such that  $ans$  does not contain the information of  $u \in [t]$  such that  $f(u) = v$ , it must be the case that  $ans$  leads to a certificate for  $x$ . It follows that the TFNP problem corresponding to  $\phi$  reduces to  $\text{rwPHP}(\text{PLS})$ .  $\square$

On the other hand, it is easy to see that  $T_2^1 + \text{dwPHP}(\text{PV})$  proves the totality of  $\text{rwPHP}(\text{PLS})$  and the proof relativizes. Let  $(f, \{I_y\}, \{g_y\})$  be an instance of  $\text{rwPHP}(\text{PLS})$ . By  $\text{dwPHP}(f)$ , there exists  $y \in [N]$  that is not in the range of  $f$ . Since  $T_2^1$  proves the totality of PLS, it also proves the existence of a solution  $ans$  for  $I_y$ . Note that since  $y$  is not in the range of  $f$ , we have in particular that  $f(g_y(ans)) \neq y$ . Hence,  $(y, ans)$  is a valid solution for this  $\text{rwPHP}(\text{PLS})$  instance.

## 4 Refuters for the Pigeonhole Principle

In this section, we study the refuter problems where the family of hard tautologies is the Pigeonhole principle  $\text{PHP}_{(n+1) \rightarrow n}$ . Our main results are the PLS-memberships of width refuter problems for (variants of)  $\text{PHP}_{(n+1) \rightarrow n}$  and the  $\text{rwPHP}(\text{PLS})$ -memberships of the size refuter problems for  $\text{PHP}_{(n+1) \rightarrow n}$ . Looking ahead, in Section 5, we will establish *universal* PLS-hardness for width refuters and  $\text{rwPHP}(\text{PLS})$ -hardness for size refuters, thereby characterizing their complexities in  $\text{TFNP}^{\text{dt}}$ .

### 4.1 Refuters for Narrow Resolution Proofs

Historically, proving size lower bounds for resolution has been challenging and considered milestones in proof complexity. The honor of the first super-polynomial size lower bounds belongs to the Pigeonhole Principle ( $\text{PHP}_{(n+1) \rightarrow n}$ ). However, in terms of width lower bound,  $\text{PHP}$  is not satisfactory:  $\text{width}(\text{PHP}_{(n+1) \rightarrow n} \vdash_{\text{Res}} \perp) = n$ , but there is already an *axiom* in  $\text{PHP}$  that has width  $n$ . Therefore, studying the complexity of finding a wide clause is uninteresting, as one of the widest clauses appears directly in the axiom and can be easily located. In what follows, we consider the width refuter problem of a *constant-width analog* of  $\text{PHP}_{(n+1) \rightarrow n}$ , namely, the nondeterministic extension  $\text{EPHP}_{(n+1) \rightarrow n}$ , defined below. However, we will come back to  $\text{PHP}_{(n+1) \rightarrow n}$  shortly after  $\text{EPHP}_{(n+1) \rightarrow n}$  and examine the refuter problem for the so-called “monotone width” of  $\text{PHP}_{(n+1) \rightarrow n}$ . The monotone width lower bound for  $\text{PHP}_{(n+1) \rightarrow n}$  will also serve as a key component in the study of size refuters for  $\text{PHP}_{(n+1) \rightarrow n}$ .

**Definition 4.1** ( $\text{EPHP}_{(n+1) \rightarrow n}$ ).  $\text{EPHP}_{(n+1) \rightarrow n}$  is the same as  $\text{PHP}_{(n+1) \rightarrow n}$  except that we replace every clause  $\bigvee_{j \in [n]} x_{ij}$  by a 3-CNF nondeterministic extension; that is, by the following  $n + 2$  clauses:

$$\bar{y}_{i,0}, (y_{i,0} \vee x_{i,1} \vee \bar{y}_{i,1}), (y_{i,1} \vee x_{i,2} \vee \bar{y}_{i,2}), \dots, (y_{i,n-1} \vee x_{i,n} \vee \bar{y}_{i,n}), y_{i,n},$$

where  $y_{i,0}, \dots, y_{i,n}$  are newly introduced variables.

Width lower bounds for  $\text{EPHP}_{(n+1) \rightarrow n}$  were proved by Ben-Sasson and Wigderson [BW01].

**Theorem 4.2** ([BW01, Theorem 4.9]). *Any resolution refutation of  $\text{EPHP}_{(n+1) \rightarrow n}$  contains a clause  $C$  with  $w(C) \geq n/3$ .*

**Theorem 4.3.**  $\text{REFUTER}(w(\text{EPHP} \vdash_{\text{Res}} \perp) < n/3)$  *is in PLS. In particular, there is a uniform decision tree reduction of block-depth 3 from the refuter problem to ITER.*

*Proof.* We will reduce it to an instance of reversed ITER. This reduction is an analog of a *constructive* version of width lower bound proofs by Beame and Pitassi [BP96], which we will use again later in the proof of [Theorem 4.8](#). We call  $(x_{i,j})$  *original* variables and  $(y_{i,j})$  *extension* variables.

We consider a set of functions over all the variables, including  $n + 1$  pigeon functions  $\{EP_i\}$  where

$$EP_i := \bar{y}_{i,0} \wedge (y_{i,0} \vee x_{i,1} \vee \bar{y}_{i,1}) \wedge (y_{i,1} \vee x_{i,2} \vee \bar{y}_{i,2}) \wedge \cdots \wedge (y_{i,n-1} \vee x_{i,n} \vee \bar{y}_{i,n}) \wedge y_{i,n}, \quad (4)$$

and  $n^2(n + 1)/2$  hole functions  $\{H_{(i,i')}^j\}$ , where  $H_{(i,i')}^j = \bar{x}_{i,j} \vee \bar{x}_{i',j}$ . It is easy to see that the semantic meaning of  $y_{i,j}$  is whether “the index of the hole that pigeon  $i$  goes into belongs to  $\{1, \dots, j\}$ ” or not.

We say an assignment  $\alpha$  over all variables (including the extension variables) is  $\ell$ -critical if  $EP_\ell(\alpha) = 0$  but all other functions are 1 under this assignment, namely  $EP_i(\alpha) = 1$  for all  $i \neq \ell$  and  $H_{(i,i')}^j(\alpha) = 1$  for all  $j \in [n], i, i' \in [n + 1]$ . If we ignore the extension variables,  $\alpha$  is essentially a complete matching between pigeons and holes, except pigeon  $i$  is not going anywhere. Given the definition of  $\ell$ -critical assignments, we define a complexity measure for a clause  $C$ , denoted by  $\text{cri}(C)$ :

$$\text{cri}(C) := |\{\ell : \exists \ell\text{-critical assignment } \alpha \text{ such that } C(\alpha) = 0\}|.$$

Note that  $\text{cri}$  has four important properties:

- (I)  $\text{cri}(\perp) = n + 1$ ;
- (II)  $\text{cri}(EP_i) = 1$  for all  $i$  and  $\text{cri}(H_{i,i'}^j) = 0$  for all  $j, i, i'$ ;
- (III)  $\text{cri}$  is subadditive with respect to resolution derivation, namely, if  $C$  is resolved from  $A$  and  $B$ , then  $\text{cri}(C) \leq \text{cri}(A) + \text{cri}(B)$ ;
- (IV) if  $C$  is obtained from a weakening of  $A$ , then  $\text{cri}(C) \leq \text{cri}(A)$ .

We first show that  $\text{cri}(\cdot)$  can be computed in polynomial time. Then we show that any clause  $C_i$  such that  $n/3 \leq \text{cri}(C_i) \leq 2n/3$  will give us a solution. The PLS-membership follows from that the standard  $1/3$ - $2/3$  trick can be implemented via a reduction to reversed ITER.

**Lemma 4.4.** *For any clause  $C$ ,  $\text{cri}(C)$  can be computed in  $\text{poly}(n)$  time.*

*Proof.* Fix any clause  $C$ . We will enumerate  $\ell$  and check the existence of  $\ell$ -critical assignments. The only part that we need to be careful is how we deal with the extension variables.

Imagine that we maintain a complete bipartite graph with  $n + 1$  nodes on the left and  $n$  nodes on the right. We will iteratively delete some edges from this graph based on the requirement of  $\ell$ -critical assignments. We will show that the existence of  $\ell$ -critical assignment can be reduced to the existence of a perfect matching of the final graph.

For an  $\ell$ -critical assignment,  $EP_\ell$  needs to be 0, and all other functions need to be 1, then we have that  $x_{\ell,j}$  needs to be 0 for all  $j$  (which means pigeon  $\ell$  cannot go into any hole). This is because if pigeon  $\ell$  were matched with some hole, some other pigeons would have no place. So we delete edges between  $(\ell, j)$  for all  $j$ .

For some  $\alpha$  being an  $\ell$ -critical assignment, we have  $C(\alpha) = 0$ . This means all literals that appeared in  $C$  are fixed to be 0 in the search of  $\alpha$ . If  $C$  contains a literal  $\bar{x}_{\ell,j}$  for some  $j$ , then we directly conclude that  $\ell$ -critical assignment does not exist (due to the argument above).

Now assume that  $C$  does not contain any literal  $\bar{x}_{\ell,j}$ . For every literal  $x_{i,j}$  in  $C$ , in order to falsify  $C$ ,  $x_{i,j}$  is going to be 0, so we delete the edge  $(i, j)$ . For every literal  $\bar{x}_{i,j}$  in  $C$ ,  $x_{i,j}$  is going to be 1, meaning that pigeon  $i$  is going to be matched with hole  $j$ , so we delete the edge  $(i, j')$  for every  $j' \neq j$  and  $(i', j)$  for every  $i' \neq i$ .

For every literal  $y_{i,j}$  in  $C$ ,  $y_{i,j}$  is going to be 0, so we delete edges  $(i, j')$  for all  $j' \leq j$ . For every literal  $\bar{y}_{i,j}$  in  $C$ ,  $y_{i,j}$  is going to be 1, so we delete edges  $(i, j')$  for all  $j' > j$ .

We can conclude that there is an  $\ell$ -critical assignment if and only if the remaining bipartite graph has a perfect matching between pigeons  $[n+1] \setminus \{\ell\}$  and holes  $[n]$ .  $\diamond$

**Reduction to ITER:** The reversed ITER instance is defined by the following function  $S : [L] \rightarrow [L]$ . For every  $i \in [L]$ :

- If  $\text{cri}(C_i) < \frac{2n}{3}$ , then  $S(i) = i$ .
- Otherwise, if  $C_i$  is a weakening from  $C_j$ , then let  $S(i) = j$ .
- Finally, if  $C_i$  is resolved from  $C_j$  and  $C_k$ : If  $\text{cri}(C_j) \geq \text{cri}(C_k)$ , then  $S(i) = j$ ; otherwise  $S(i) = k$ .

It is easy to see that this reduction can be implemented in block-depth 3: for example, if  $C_i$  is resolved from  $C_j$  and  $C_k$ , then one only needs to read the  $i$ -th,  $j$ -th, and  $k$ -th node in the resolution refutation.

Note that when we find any solution  $i$  of this reversed ITER instance, it satisfies  $S(i) < i$  and  $S(S(i)) = S(i)$ . This means  $\text{cri}(C_i) \geq 2n/3$  but  $\text{cri}(C_{S(i)}) < 2n/3$ . Thus we have  $\text{cri}(C_{S(i)}) \in [n/3, 2n/3]$ .

Now it remains for us to show that any  $C$  such that  $n/3 \leq \text{cri}(C) \leq 2n/3$  has width at least  $n/3$ . For the sake of contraction, assume that  $\text{width}(C) < n/3$ . This implies that for at most  $n/3$  pigeon  $i$ ,  $C$  has some variable related to  $i$ , namely, some  $x_{i,j}$  or  $y_{i,j}$ . Since  $\text{cri}(C) \geq n/3$ , we know that there exists  $\ell$  such that there is an  $\ell$ -critical assignment  $\alpha$  for  $C$  but  $C$  has no variables of the form  $x_{\ell,j}$  or  $y_{\ell,j}$ .

On the other hand, since  $\text{width}(C) < n/3$  and  $\text{cri}(C) \leq 2n/3$ , we know that there is another index  $\ell'$  such that  $\ell'$  is not critical to  $C$  and  $C$  has no variables of the form  $x_{\ell',j}$  or  $y_{\ell',j}$ . Let  $k$  be the hole that is matched with  $\ell'$  in  $\alpha$ . Consider the following assignment  $\alpha'$ : we start from  $\alpha' := \alpha$ , flip  $x_{\ell,k}$  from 0 to 1, and flip  $x_{\ell',k}$  from 1 to 0. We further flip all  $y_{\ell,j}$  and  $y_{\ell',j}$  correspondingly. We have that  $EP_{\ell}(\alpha') = 1$  and  $EP_{\ell'}(\alpha') = 0$ . Since  $C$  doesn't contain any variables related to pigeons  $\ell$  and  $\ell'$ ,  $C(\alpha')$  is still 0. This constructs a witness that  $\ell'$  is critical to  $C$ , a contradiction. This finishes the proof.  $\square$

It is easy to see that the above reduction is actually a formalization of the width lower bound in  $\mathsf{T}_2^1(\alpha)$ . Let  $n \in \text{Log}$ ,  $\alpha(\cdot, \cdot)$  be an oracle that encodes a purported resolution proof of  $\text{EPHP}_{(n+1) \rightarrow n}$  with width less than  $n/3$ , where the second input is a parameter (non-uniformity); that is, each  $z$  corresponds to a resolution proof  $\pi_z$  and  $\alpha(i, z)$  is the  $i$ -th bit of  $\pi_z$ . We can syntactically enforce every clause in the proof to have width at most  $n/3$ . Let  $\text{mistake}_{\text{EPHP}}^{\text{w}}(n, \alpha, z, i)$  denote the PV( $\alpha$ ) predicate stating that the  $i$ -th derivation step in  $\pi_z$  is invalid (the superscript “w” stands for “width”). Then we have:

**Theorem 4.5.**

$$\mathsf{T}_2^1(\alpha) \vdash \forall n \in \text{Log} \forall z \exists i \text{mistake}_{\text{EPHP}}^{\text{w}}(n, \alpha, z, i).$$

*Proof Sketch.* Reason in  $\mathsf{T}_2^1(\alpha)$ . Assuming  $\forall i \neg \text{mistake}_{\text{EPHP}}^{\text{w}}(n, \alpha, z, i)$ , we will derive a contradiction. A minor technical issue is that we need a PV-definition of the function  $\text{cri}$  such that the properties (I)-(IV) are true, and that any clause  $C$  with  $n/3 \leq \text{cri}(C) \leq 2n/3$  has width at least  $n/3$ . This follows from the formalization of bipartite matching algorithms in PV [LC11].<sup>23</sup>

Let  $C_1, C_2, \dots, C_L$  denote the purported resolution refutation encoded by  $\alpha(\cdot, z)$ . Using  $\Sigma_1^b(\alpha)$ -MIN (which is available in  $\mathsf{T}_2^1(\alpha)$ ), there is a smallest integer  $i$  such that  $\text{cri}(C_i) \geq n/3$ . By (II),  $C_i$  cannot

<sup>23</sup>This annoying detail would disappear if we consider the universal variant  $\forall \mathsf{T}_2^1(\alpha)$  since these properties are indeed true universal sentences in the standard PV model  $\mathbb{N}$  and thus are included in the axioms of  $\forall \mathsf{T}_2^1(\alpha)$ . As pointed out in [Mül21], it is the provability in *universal* variants of relativized bounded arithmetic that captures reducibility among type-2 TFNP problems.

be an axiom of  $\text{EPHP}_{(n+1) \rightarrow n}$ . By (IV),  $C_i$  cannot be a weakening of any clause  $C_j$  ( $j < i$ ), as this would contradict the minimality of  $i$ . Hence  $C_i$  is resolved from  $C_j$  and  $C_k$  for some  $j, k < i$ . By (III),  $\text{cri}(C_i) \leq \text{cri}(C_j) + \text{cri}(C_k) \leq 2n/3$ . Hence the width of  $C_i$  is at least  $n/3$ , a contradiction.  $\square$

**Monotonized resolution and its width refuter.** The first exponential size lower bound for resolution was proven by Haken [Hak85] for the pigeonhole principle PHP. Haken used the so-called ‘‘bottleneck counting’’ argument and the proof was quite involved. A much simpler proof was found by Beame and Pitassi [BP96], where one of the crucial ingredients is the following lemma.

**Lemma 4.6** ([BP96]). *Any resolution refutation of  $\text{PHP}_{(n+1) \rightarrow n}$  contains a clause  $C$  with  $w(\text{mono}(C)) \geq 2n^2/9$ .*

Here, for a clause  $C$ ,  $\text{mono}(C)$  is the ‘‘monotonized’’ version of  $C$ , which is obtained from  $C$  by replacing every negated variable  $\bar{x}_{ij}$  with the set of variables  $x_{i'j}$  for all  $i' \neq i$ .<sup>24</sup> Given Lemma 4.6, we could define another variant of the width refuter problem that concerns the original PHP tautologies with no extension variables: we wish to find a clause  $C$  such that  $\text{mono}(C)$  has a large width, in particular,  $w(\text{mono}(C)) \geq 2n^2/9$ . The problem is denoted as  $\text{REFUTER}(w_{\text{mono}}(\text{PHP} \vdash_{\text{Res}} \perp) < 2n^2/9)$  and is defined as follows.

**Definition 4.7** ( $\text{REFUTER}(w_{\text{mono}}(\text{PHP} \vdash_{\text{Res}} \perp))$ ). Consider the tautology  $\text{PHP}_{(n+1) \rightarrow n}$  and let  $w := 2n^2/9$ . The input instance  $\Pi$  is a purported resolution proof for  $\text{PHP}_{(n+1) \rightarrow n}$  that consists of clauses  $C_{-k}, \dots, C_{-1}, C_0, \dots, C_{L-1}$ , where the first  $k := n + 1 + n^2(n + 1)/2$  clauses are axioms from  $\text{PHP}_{(n+1) \rightarrow n}$  and the last clause  $C_L = \perp$ .

A solution of the given instance is one of the following:

- an index  $i \in [L]$  such that  $\text{mono}(C_i)$  has at least  $w$  literals;<sup>25</sup> or
- an index  $i$  such that  $C_i$  is an invalid derivation.

The width refuter problem of monotized resolution proof may not seem natural in the first place. However, converting resolution proof into monotized resolution proof is an elegant ingredient in the *size lower bounds* of PHP in the proof of Beame and Pitassi [BP96]. Ultimately, our main motivation is for the size refuter of PHP, and the PLS-membership of  $\text{REFUTER}(w_{\text{mono}}(\text{PHP} \vdash_{\text{Res}} \perp) < 2n^2/9)$  is a key step of showing the  $\text{rwPHP}(\text{PLS})$ -membership of the size refuter ( $\text{REFUTER}(s(\text{PHP} \vdash_{\text{Res}} \perp) < 1.01^n)$ ). Indeed, the  $\text{rwPHP}(\text{PLS})$ -membership in Section 4.2 uses the following theorem as a black box:

**Theorem 4.8.**  $\text{REFUTER}(w_{\text{mono}}(\text{PHP}_{(n+1) \rightarrow n} \vdash_{\text{Res}} \perp) < 2n^2/9)$  is in PLS. In particular, there is a uniform decision tree reduction of block-width 3 from this refuter problem to ITER.

The proof is in fact simpler than that of Theorem 4.3, and is a proper constructive translation of the proof by Beame and Pitassi [BP96].

*Proof of Theorem 4.8.* We call an assignment  $\alpha$  to be  $\ell$ -critical if  $\alpha$  is a perfect matching between pigeons and holes, except pigeon  $\ell$  is not going anywhere. Formally, for every  $i \neq \ell$  we have  $x_{i1} \vee \dots \vee x_{in}$ , and for every  $i, i', j$  we have  $\bar{x}_{ij} \vee \bar{x}_{i'j}$  under  $\alpha$ .

Given the definition of  $\ell$ -critical assignments, we define  $\text{cri}(\text{mono}(C))$  as follows:

$$\text{cri}(\text{mono}(C)) := |\{\ell : \exists \ell\text{-critical assignment } \alpha \text{ that falsifies } \text{mono}(C)\}|.$$

Again,  $\text{cri}$  has the following four important properties:

<sup>24</sup>The notion of  $\text{mono}(C)$  is tailored to the hard tautology  $\text{PHP}_{(n+1) \rightarrow n}$ . The proof of [BP96] only considers assignments  $x \in \{0, 1\}^{(n+1)n}$  that defines a bijective mapping from  $n$  of the pigeons to all  $n$  holes; it is easy to see that every clause  $C$  is equivalent to  $\text{mono}(C)$  w.r.t. such ‘‘critical’’ assignments  $x$ .

<sup>25</sup>Here, unlike the formalization of width lower bounds in Definition 3.6, it is unclear how to syntactically enforce that the *monotonized version* of every clause has width  $< w$ . Therefore we include the clauses with large monotone width as solutions.

- $\text{cri}(\perp) = n + 1$ ;
- $\text{cri}(\text{mono}(C)) \leq 1$  for all axioms  $C$  of  $\text{PHP}_{(n+1) \rightarrow n}$ ;
- $\text{cri}$  is subadditive with respect to resolution derivation, namely, if  $C$  is derived from  $A$  and  $B$ , then  $\text{cri}(\text{mono}(C)) \leq \text{cri}(\text{mono}(A)) + \text{cri}(\text{mono}(B))$ .
- If  $C$  is obtained from a weakening of  $A$ , then  $\text{cri}(\text{mono}(C)) \leq \text{cri}(\text{mono}(A))$ .

This time, since we are only concerned with the monotonized version of a clause  $C$  and there are no extension variables, it is easier to show that  $\text{cri}(\text{mono}(C))$  can be computed in polynomial time.

**Claim 4.9.** *For any clause  $C$ ,  $\text{cri}(\text{mono}(C))$  can be computed in polynomial time.*

*Proof.* Fix a pigeon  $\ell$  and we want to check if there is an  $\ell$ -critical assignment for  $C$ . We maintain a complete bipartite graph and delete all edges between  $(\ell, j)$  for all hole  $j$ . If a variable  $x_{ij}$  appears in  $C$ , we delete the edge  $(i, j)$ . Then  $\ell$ -critical assignment exists if and only if there is a perfect matching between pigeons  $[n + 1] \setminus \{\ell\}$  and holes  $[n]$ .  $\diamond$

**Claim 4.10.** *For any clause  $C$ , the width of  $\text{mono}(C)$  is at least  $\text{cri}(\text{mono}(C)) \cdot (n - \text{cri}(\text{mono}(C)))$ .*

*Proof.* Let  $D := \text{mono}(C)$ . Let  $\text{CriP}(D)$  be the critical pigeons to  $D$ , i.e., the set of pigeons  $\ell \in [n + 1]$  such that there exists an  $\ell$ -critical assignment falsifying  $D$ . Then  $\text{cri}(\text{mono}(C)) = |\text{CriP}(D)|$ . Let  $u_1 \in \text{CriP}(D)$  and  $u_2 \notin \text{CriP}(D)$ . Since  $u_1 \in \text{CriP}(D)$ , there is a  $u_1$ -critical assignment  $\alpha$  that falsifies  $D$ . Suppose that  $u_2$  is mapped to the hole  $v_2$  in the assignment  $\alpha$ . Let  $\beta$  denote the assignment obtained from  $\alpha$  by mapping  $u_1$  into  $v_2$  and not mapping  $u_2$  anywhere. Then  $\beta$  is a  $u_2$ -critical assignment. Since  $u_2 \notin \text{CriP}(D)$ ,  $\beta$  satisfies  $D$ . However, there is only one variable that appears positively in  $\beta$  but negatively in  $\alpha$ : namely,  $x_{u_1, v_2}$ . Since  $D$  is monotone, the literal  $x_{u_1, v_2}$  appears in  $D$ . Repeating this argument for every  $u_1 \in \text{CriP}(D)$  and  $u_2 \notin \text{CriP}(D)$ , we can see that the width of  $D$  is at least  $\text{cri}(D) \cdot (n - \text{cri}(D))$ .  $\diamond$

Given the lemma above, it remains for us to show that finding a clause  $C$  such that  $\text{cri}(\text{mono}(C)) \in [n/3, 2n/3]$  belongs in PLS. This can be implemented by the standard 1/3-2/3 trick in a potential function way, which is exactly the same as the argument used in the proof of [Theorem 4.3](#).  $\square$

## 4.2 Refuters for Short Resolution Proofs

In this subsection, we investigate [Problem 1.1](#), i.e., the refuter for the following classic resolution *size* lower bound:

**Theorem 4.11** ([\[Hak85\]](#)). *Any resolution refutation of  $\text{PHP}_{(n+1) \rightarrow n}$  requires at least  $2^{\Omega(n)}$  clauses.*

We show that this refuter problem is in  $\text{rwPHP}(\text{PLS})$ . As mentioned in [Section 1.3.2](#), this is done by carefully following the proofs of [Theorem 4.11](#). We follow the simplified proof by Beame and Pitassi [\[BP96\]](#), which consists of two steps: a (monotone) width lower bound and a random restriction argument. As the required width lower bound was already studied in [Theorem 4.8](#), we focus on the random restriction argument here.

Let  $X := [n + 1]$  denote the set of pigeons and  $Y := [n]$  denote the set of holes. Recall that the *monotone* version of a clause  $C$ , denoted as  $\text{mono}(C)$ , is obtained from  $C$  by replacing every negated variable  $\bar{x}_{ij}$  with the set of variables  $x_{ij'}$  for all  $j' \neq j$ .

**Definition 4.12.** Let  $t < n$  be a parameter,  $\pi : X \rightarrow Y$  be a *size- $t$  matching*, i.e., a partial injective function with  $|\text{Domain}(\pi)| = t$ . This matching induces a restriction that sets:

- the variable  $x_{u, \pi(u)}$  to be 1, for every  $u \in \text{Domain}(\pi)$ ;
- the variable  $x_{u, v}$  to be 0, for every  $u \in \text{Domain}(\pi)$  and  $v \in Y \setminus \{\pi(u)\}$ ; and

- the variable  $x_{u',\pi(u)}$  to be 0, for every  $u \in \text{Domain}(\pi)$  and  $u' \in X \setminus \{u\}$ .

Suppose that  $C_0, C_1, \dots, C_{L-1}$  is a resolution refutation of  $\text{PHP}_{(n+1) \rightarrow n}$ . For each clause  $C_i$ , let  $\pi(C_i)$  denote the sub-clause of  $C_i$  under the above restriction. That is, if  $\pi$  sets some variable in  $C_i$  to 1 then  $\pi(C_i) = 1$ ; otherwise  $\pi(C_i)$  is obtained from  $C_i$  by removing every variable set to 0 by  $\pi$ . Note that the above restriction transforms the unsatisfiable CNF  $\text{PHP}_{(n+1) \rightarrow n}$  into the unsatisfiable CNF  $\text{PHP}_{(n-t+1) \rightarrow (n-t)}$ . Then,  $\pi(C_0), \pi(C_1), \dots, \pi(C_{L-1})$  is a resolution refutation for  $\text{PHP}_{(n-t+1) \rightarrow (n-t)}$ . Now we claim that the width of the resolution refutation becomes small after being restricted by a random matching  $\pi$ .

**Claim 4.13.** *If a size- $t$  matching  $\pi$  is chosen uniformly at random over all possible size- $t$  matchings, then with probability at least  $1/2$ , it holds that for every clause  $i \in [L]$ ,  $w(\text{mono}(\pi(C_i))) < W$ , where  $W := (n+1)^2(1 - (1/2L)^{1/t})$ .*

*Proof.* Fix  $i \in [L]$ , we show that the probability over  $\pi$  that  $w(\text{mono}(\pi(C_i))) > W$  is at most  $1/(2L)$ . The claim then follows from a union bound.

Choose the matching  $\pi$  round by round. There are  $t$  rounds, where in each round, we choose an unmatched  $u \in X$  and an unmatched  $v \in Y$  uniformly at random and match them. If, for the current partial matching  $\pi$ , we have  $w(\text{mono}(\pi(C_i))) \geq W$ , then the probability that  $x_{u,v} \in \text{mono}(\pi(C_i))$  is at least  $W/(n+1)^2$ . If this is the case, then  $\text{mono}(\pi(C_i))$  will become 1 (the always-true clause) after we set  $\pi(u) \leftarrow v$ , thus it gets “killed.” It follows that the probability that  $C_i$  never gets killed is at most  $(1 - W/(n+1)^2)^t \leq 1/(2L)$ .  $\square$

Combining [Lemma 4.6](#) and [Claim 4.13](#), we obtain the following size lower bound:

**Theorem 4.14.** *Any resolution refutation of  $\text{PHP}_{(n+1) \rightarrow n}$  requires more than  $L := 1.01^n$  clauses.*

*Proof.* Let  $t := n/10$ , then  $W = (n+1)^2(1 - (1/2L)^{1/t}) \leq \frac{2}{9}(n-t)^2$ . If there is a resolution refutation of  $\text{PHP}_{(n+1) \rightarrow n}$  of size at most  $L$ , then by [Claim 4.13](#), there is a resolution refutation of  $\text{PHP}_{(n-t+1) \rightarrow (n-t)}$  of monotone width at most  $\frac{2}{9}(n-t)^2$ . This contradicts [Lemma 4.6](#).  $\square$

To derive an upper bound for the complexity of refuter that corresponds to [Theorem 4.14](#), we need a *constructive* version of [Claim 4.13](#). We start by setting up an encoding for the partial matchings and random restrictions that will make it easier to describe our reductions.

A size- $t$  matching can be described by an edge-sequence  $(u_0, v_0), (u_1, v_1), \dots, (u_{t-1}, v_{t-1})$ , where for each  $j \in [t]$ ,  $u_j \in [n-j+1]$  and  $v_j \in [n-j]$ . The first edge in this matching connects the  $u_0$ -th node in  $X$  and the  $v_0$ -th node in  $Y$  (the first node is the 0-th), the second edge connects the  $u_1$ -th unused node in  $X$  (i.e.,  $u_1$ -th node in  $X \setminus \{u_1\}$ ) and the  $v_1$ -th unused node in  $Y$ , and so on.<sup>26</sup> The space of all possible edge-sequences is denoted by

$$\mathcal{SEQ} := ([n+1] \times [n]) \times ([n] \times [n-1]) \times \dots \times ([n-t+2] \times [n-t+1]).$$

On the other hand, fix a clause  $C_i$  such that  $w(\text{mono}(C_i)) \geq W$ . Say an edge-sequence  $s$  is *bad* for  $C_i$  if  $w(\text{mono}(s(C_i))) \geq W$ , where  $\pi_s(C_i)$  is the restriction of  $C_i$  under the matching corresponding to  $\pi_s$ . As we argued in [Claim 4.13](#), the number of bad edge-sequences for each  $C_i$  is small; we set up another encoding to justify this fact. Any bad edge-sequence can be encoded as a sequence  $(e_0, e_1, \dots, e_{t-1})$ ,<sup>27</sup> where for each  $j \in [t]$ ,  $e_j \in [(n-j+1)(n-j) - W]$ . The first edge  $(u_0, v_0)$  in this matching is the  $e_0$ -th edge, in the lexicographical order, that is not a literal in  $\text{mono}(C_i)$ ; the second edge  $(u_1, v_1)$  is the  $e_1$ -th edge that still can be chosen (we cannot choose any edge touching either  $u_0$  or  $v_0$ ) and is not a literal in the current  $\text{mono}(\pi_s(C_i))$ ; and so on. If  $s$  is bad, then  $w(\text{mono}(\pi_s(C_i)))$  never goes below  $W$ , therefore at

<sup>26</sup>Note that the edges are *ordered*, hence each matching corresponds to  $t!$  different edge-sequences. In what follows, we will talk about edge-sequences instead of matchings.

<sup>27</sup>To avoid confusion, we use “edge-sequence” to denote elements in  $\mathcal{SEQ}$  and “sequence” to denote elements in  $\mathcal{BAD}$ .

the  $j$ -th stage, there are at most  $(n - j + 1)(n - j) - W$  possible edges to choose. Hence, the space of all possible sequences encoding bad edge-sequences is:

$$\mathcal{BAD} = [(n + 1)n - W] \times [n(n - 1) - W] \times \cdots \times [(n - t + 2)(n - t + 1) - W].$$

The following calculation corresponds to a *union bound* over all  $C_i$  that the number of bad edge-sequences is small:

$$\begin{aligned} \frac{|\mathcal{BAD}| \cdot L}{|\mathcal{SEQ}|} &= L \cdot \prod_{j=0}^{t-1} \frac{(n + 1 - j)(n - j) - W}{(n + 1 - j)(n - j)} \\ &\leq L \cdot (1 - W/(n + 1)^2)^t \\ &\leq 1/2. \end{aligned} \tag{5}$$

Fix a clause  $C_i$ . For every sequence  $b \in \mathcal{BAD}$ , let  $\text{seq}(C_i, b) \in \mathcal{SEQ}$  denote the bad edge-sequence for  $C_i$  corresponding to  $b$ ; if  $\text{seq}(C_i, b)$  does not exist<sup>28</sup>, then we denote  $\text{seq}(C_i, b) = \perp$ . Conversely, any  $s \in \mathcal{SEQ}$  is either bad for  $C_i$  or not; if  $s$  is bad for  $C_i$ , then denote  $b := \text{bad}(C_i, s)$  as the sequence  $b \in \mathcal{BAD}$  corresponding to  $s$ ; otherwise we say  $\text{bad}(C_i, s) := \perp$ . We need the fact that:

**Fact 4.15.** *Let  $s \in \mathcal{SEQ}$  be bad for the clause  $C_i$ , then  $\text{seq}(C_i, \text{bad}(C_i, s)) = s$ .*

Now we are ready to establish the rwPHP(PLS) upper bound for the refuter of [Theorem 4.14](#).

**Theorem 4.16.** *There is a uniform decision tree reduction of block-depth 3 from  $\text{REFUTER}(\text{Res}(\text{PHP}) > 1.01^n)$  to  $\text{rwPHP}(\text{PLS})$ .*

*That is, there is a uniform decision tree reduction of block-depth 3 such that the following holds:*

- *given a resolution refutation  $\Pi = (C_0, C_1, \dots, C_{L-1})$  for  $\text{PHP}_{(n+1) \rightarrow n}$ , where  $L \leq 1.01^n$ , the reduction computes an instance  $(f, \{I_y\}, \{g_y\})$  of  $\text{rwPHP}(\text{PLS})$ ;*
- *given any valid answer for  $(f, \{I_y\}, \{g_y\})$ , one can compute an invalid derivation  $C_i \in \Pi$  in  $\text{poly}(n)$  time.*

*Proof.* Let  $M := |\mathcal{BAD}| \cdot L$  and  $N := |\mathcal{SEQ}|$ , then from [Equation 5](#), we have  $M \leq N/2$ . We will identify numbers in  $[M]$  with pairs  $(i, b)$  where  $i \in [L]$  and  $b \in \mathcal{BAD}$ , and identify numbers in  $[N]$  with edge-sequences in  $\mathcal{SEQ}$ . The instance  $(f, \{I_y\}, \{g_y\})$  is defined as follows:

( $f$ ) For every  $x \in [M]$ , we interpret  $x$  as a pair  $(i, b)$  where  $i \in [L]$  and  $b \in \mathcal{BAD}$ . If  $\text{seq}(C_i, b) \neq \perp$ , then we let  $f(x) := \text{seq}(C_i, b)$ ; otherwise let  $f(x) := 0$  (the choice 0 is arbitrary).

( $I_y$ ) Fix  $y \in [N] = \mathcal{SEQ}$ . The edge-sequence  $y$  defines a size- $t$  partial matching  $\pi_y$ , which induces a resolution refutation  $\Pi|_y = (\pi_y(C_0), \pi_y(C_1), \dots, \pi_y(C_{L-1}))$  of  $\text{PHP}_{(n-t+1) \rightarrow (n-t)}$ . We treat  $\Pi|_y$  as a purported resolution refutation with monotone width  $< W$ ; by [Theorem 4.8](#), the problem of finding an invalid derivation in  $\Pi|_y$  reduces to ITER via a decision tree of block-width 3. Let  $I_y$  be the ITER instance obtained by this reduction.

( $g_y$ ) Fix  $y \in [N] = \mathcal{SEQ}$  and an answer  $ans$  of the ITER instance  $I_y$ . Given  $ans$ , we can compute a clause that is either invalid or too fat; we then compute  $g_y(ans)$  from this clause.

More precisely, we can compute an integer  $i \in [L]$  such that either  $\text{width}(\text{mono}(\pi_y(C_i))) \geq W$ , or  $\pi_y(C_i)$  corresponds to an invalid derivation in  $\Pi|_y$ . In the second case,  $C_i$  is also an invalid derivation in  $\Pi$ , thus we can set  $g_y(ans)$  to be an arbitrary value (say 0). In the first case, we can set  $g_y(ans) := (i, \text{bad}(C_i, y))$ .

<sup>28</sup>This may happen when, for example,  $w(\text{mono}(C_i))$  is much larger than  $W$  and  $b_0 > (n + 1)n - w(\text{mono}(C_i))$ .

The block-depth of the decision trees computing  $f$ ,  $\{I_y\}$ , and  $\{g_y\}$  are 1, 3, and 2 respectively. Clearly, the decision trees are uniform.

Finally, let  $(y, ans)$  be a valid solution for the rwPHP(PLS) instance  $(f, \{I_y\}, \{g_y\})$  (i.e.,  $f(g_y(ans)) \neq y$ ). The edge-sequence  $y \in \mathcal{SEQ}$  corresponds to a size- $t$  partial matching  $\pi_y$  and from  $ans$  we can read off a clause  $i \in [L]$  such that either (1)  $\text{width}(\text{mono}(\pi_y(C_i))) \geq W$  or (2)  $\pi_y(C_i)$  is an invalid derivation in  $\Pi|_y$ . If (1) holds, then

$$f(g_y(ans)) = f(i, \text{bad}(C_i, y)) = \text{seq}(C_i, \text{bad}(C_i, y)) = y$$

by [Fact 4.15](#). This contradicts  $(y, ans)$  being a valid solution for rwPHP(PLS). It follows that (2) happens and we have found an invalid derivation (namely  $C_i$ ) in  $\Pi$ .  $\square$

As mentioned in [Section 1.3.2](#), the above arguments are essentially a formalization of Haken's lower bounds in the theory  $\text{T}_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha))$ ; the decision tree reduction in [Theorem 4.16](#) follows from the witnessing theorem for  $\text{T}_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha))$  (see [Section 3.5.1](#)). In what follows, we make this formalization explicit:

**Theorem 4.17.** *Let  $n \in \text{Log}$ ,  $\alpha(\cdot, \cdot)$  be an oracle that encodes a purported length- $1.01^n$  resolution proof of  $\text{PHP}_{(n+1) \rightarrow n}$ , where the second input is a parameter. Let  $\text{mistake}_{\text{PHP}}(n, \alpha, z, i)$  denote the  $\text{PV}(\alpha)$  predicate stating that the  $i$ -th derivation step in  $\alpha(\cdot, z)$  is invalid. Then*

$$\text{T}_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha)) \vdash \forall n \in \text{Log} \forall z \exists i \in [1.01^n] \text{mistake}_{\text{PHP}}(n, \alpha, z, i).$$

*Proof Sketch.* Reason in  $\text{T}_2^1(\alpha) + \text{dwPHP}(\text{PV}(\alpha))$ ; assuming  $\forall i \in [1.01^n] \neg \text{mistake}_{\text{PHP}}(n, \alpha, z, i)$ , we will derive a contradiction. We still use  $C_i$  to denote the  $i$ -th clause of the resolution proof, noticing that given  $i$  and  $z$ ,  $C_i$  can be computed by a  $\text{PV}(\alpha)$  function. We also use our previous notation such as  $\mathcal{BAD}$  and  $\mathcal{SEQ}$ , and previous parameters  $t := n/10$ ,  $L := 1.01^n$ , and  $W := (n+1)^2(1 - (1/2L)^{1/t}) \leq \frac{2}{9}(n-t)^2$ .

First, we use  $\text{dwPHP}(\text{PV}(\alpha))$  to select a good random restriction under which each  $C_i$  becomes a small-width clause. This random restriction will be encoded as an edge-sequence  $s \in \mathcal{SEQ}$ . Consider the function

$$\overline{\text{bad}}_z(i, b) := \text{bad}(C_i, b),$$

where  $b \in \mathcal{BAD}$  is any sequence encoding a bad edge-sequence. Clearly,  $\overline{\text{bad}}_z$  is a function symbol in  $\text{PV}(\alpha)$ . By  $\text{dwPHP}(\text{PV}(\alpha))$ , there is an edge-sequence  $s \in \mathcal{SEQ}$  such that for every  $i \in [L]$  and  $b \in \mathcal{BAD}$ ,  $\overline{\text{bad}}_z(i, b) \neq s$ .

Next, we apply  $s$  to each clause  $C_i$ ; denote  $\pi_s(C_i)$  the restriction of  $C_i$  under the matching corresponding to  $s$ . By our choice of  $s$ , for every  $i \in [L]$ , we have  $w(\text{mono}(\pi_s(C_i))) < W \leq \frac{2}{9}(n-t)^2$ . By [Claim 4.10](#), we have  $\text{cri}(\text{mono}(\pi_s(C_i))) > \frac{2(n-t)}{3}$  or  $\text{cri}(\text{mono}(\pi_s(C_i))) < \frac{n-t}{3}$  for every  $i \in [L]$ .

Then we invoke [Lemma 4.6](#) to show that the sequence  $\pi_s(C_0), \pi_s(C_1), \dots, \pi_s(C_{L-1})$  is not a valid resolution proof for  $\text{PHP}_{(n-t+1) \rightarrow (n-t)}$ . Note that this is the step where we use the power of  $\text{T}_2^1(\alpha)$ . Since  $\text{cri}(\text{mono}(\pi_s(C_{L-1}))) = \text{cri}(\perp) = n+1$ , by  $\Sigma_1^b(\alpha)$ -MIN (which is available in  $\text{T}_2^1(\alpha)$ ), there is a smallest integer  $i \leq L-1$  such that  $\text{cri}(\text{mono}(\pi_s(C_i))) > \frac{2(n-t)}{3}$ .

- If  $\pi_s(C_i)$  is an axiom, then  $\text{cri}(\text{mono}(\pi_s(C_i))) \leq 1$ , which is a contradiction.
- If  $\pi_s(C_i)$  is a weakening of  $\pi_s(C_j)$  where  $j < i$ , then  $\text{cri}(\text{mono}(\pi_s(C_j))) \leq \text{cri}(\text{mono}(\pi_s(C_i)))$ , contradicting the minimality of  $i$ .
- If  $\pi_s(C_i)$  is a resolution of  $\pi_s(C_j)$  and  $\pi_s(C_k)$  where  $j, k < i$ , then

$$\text{cri}(\text{mono}(\pi_s(C_i))) \leq \text{cri}(\text{mono}(\pi_s(C_j))) + \text{cri}(\text{mono}(\pi_s(C_k))).$$

However, this means either  $\text{cri}(\text{mono}(\pi_s(C_j)))$  or  $\text{cri}(\text{mono}(\pi_s(C_k)))$  is at least  $\frac{n-t}{3}$ , a contradiction.

Now, since  $\pi_s(C_0), \pi_s(C_1), \dots, \pi_s(C_{L-1})$  is not a valid resolution proof for  $\text{PHP}_{(n-t+1) \rightarrow (n-t)}$ , we have that  $C_0, C_1, \dots, C_{L-1}$  is not a valid resolution proof of  $\text{PHP}_{(n+1) \rightarrow n}$  either. This finishes the proof.  $\square$

## 5 Hardness of Refuting Resolution Proofs

In this section, we provide two hardness results for refuter problems: namely, the PLS-hardness for resolution width refuters ([Theorem 5.1](#)) and the rwPHP(PLS)-hardness for resolution size refuters ([Theorem 5.4](#)). Notably, our hardness results hold for any family of hard tautologies *as long as the lower bounds are true*. This means that “PLS-reasoning” is necessary for proving *any* non-trivial resolution width lower bounds and “rwPHP(PLS)-reasoning” is necessary for proving *any* non-trivial resolution size lower bounds. Of course, this also implies that PLS and rwPHP(PLS) are the *best possible* upper bounds for the refuter problems for *any* non-trivial unsatisfiable family of CNFs (as what we obtained for  $\text{PHP}_{(n+1) \rightarrow n}$  and more upper bounds in [Section 6](#)).

### 5.1 Hardness of Refuting Narrow Resolution Proofs

We show that the refuter problems for any *true* resolution width lower bounds are PLS-hard. In fact, this holds even for unsatisfiable CNF families that only contain a *single* CNF of *constant* size:

**Theorem 5.1.** *Let  $F$  be any unsatisfiable CNF with a non-trivial resolution width lower bound, i.e.,  $w(F \vdash_{\text{Res}} \perp) > \text{width}(F)$ . Let  $\mathcal{F} := \{F\}$  and  $w_0 := \text{width}(F)$ . Then there is a (uniform) decision-tree reduction of block-depth 2 from ITER to  $\text{REFUTER}(w(F \vdash_{\text{Res}} \perp) \leq w_0)$ .*

*Proof.* We show a straightforward reduction from the reversed ITER to the width refuter.

Note that  $F$  is a fixed CNF so it can be seen as constant size. Hence we can check in constant time that  $F$  is unsatisfiable and  $w(F \vdash_{\text{Res}} \perp) > \text{width}(F)$ .

Let  $S : [L] \rightarrow [L]$  such that  $S(L-1) < L-1$  be any instance of reversed ITER. We will construct a purported resolution refutation  $\Pi$  for  $F$  such that any invalid derivation in  $\Pi$  corresponds to an answer for  $S$ . Let  $k$  be the number of axioms in  $F$ . The resolution refutation  $\Pi$  consists of nodes  $C_{-k}, \dots, C_{-1}, C_0, \dots, C_{L-1}$ , where  $C_{-k}, \dots, C_{-1}$  are the axioms of  $F$ .

For every  $i$  such that  $S(i) = i$ , we let  $C_i = C_{-k}$  and define  $C_i$  to be a weakening from  $C_{-k}$ . This is a valid derivation (and  $C_i$  will not be used anymore). The clauses written in all other nodes in  $C_1, \dots, C_{L-1}$  will be  $\perp$ . The weakening rules applied among these nodes will encode the successor pointer  $S$ :

- For every solution  $i$  of the reversed ITER instance (i.e.,  $i$  such that either  $S(i) > i$  or  $(S(i) < i$  and  $S(S(i)) = S(i))$ ), the weakening rule applied for  $C_i$  will be *invalid*. More specifically, let  $C_i$  be a weakening from  $C_{-k}$ , then  $C_i$  becomes a solution of the  $\text{REFUTER}(w(F \vdash_{\text{Res}} \perp) \leq w_0)$  instance.
- For every  $i$  such that  $S(i) < i$  and  $S(S(i)) < S(i)$ , we let  $C_i$  be a weakening from  $C_{S(i)}$ . Since both  $C_i$  and  $C_{S(i)}$  are  $\perp$ , this is a valid derivation.

This finishes the construction, and the correctness follows from the following two facts immediately: (1) there are no nodes whose width is larger than  $\text{width}(F)$ ; (2) a resolution derivation is invalid if and only if it is a solution of the given reversed ITER instance. The block-depth of our reduction is 2, as we only need to query  $S(i)$  and  $S(S(i))$ .  $\square$

Note that the reduction in [Theorem 5.1](#) also works for a family of CNFs  $\{F_n\}_{n \in \mathbb{N}}$  with non-trivial width lower bound. Therefore, combined with the PLS-membership results in [Section 4.1](#), we obtain:

**Theorem 5.2.**  $\text{REFUTER}(w(\text{EPHP} \vdash_{\text{Res}} \perp) < n/3)$  is PLS-complete.

Note that every clause generated in the reduction in [Theorem 5.1](#) has monotone width  $O(n)$ . Hence the same proof also shows the PLS-hardness of *monotone* width refuters (as in [Theorem 4.8](#)):

**Theorem 5.3.**  $\text{REFUTER}(w_{\text{mono}}(\text{PHP}_{(n+1) \rightarrow n} \vdash_{\text{Res}} \perp) < 2n^2/9)$  is PLS-complete.

## 5.2 Hardness of Refuting Short Resolution Proofs

In this section, we show the rwPHP(PLS)-hardness of refuters for resolution *size* lower bounds. In particular, for any family of unsatisfiable CNF formulas  $\{F_n\}_{n \in \mathbb{N}}$  that requires resolution size  $> s_F(n)$ , if  $s_F(n)$  is not too small, then rwPHP(PLS) reduces to the problem  $\text{REFUTER}(s(F_n \vdash_{\text{Res}} \perp) \leq s_F(n))$ . This result and our rwPHP(PLS) upper bounds (Theorems 4.16, 6.4, 6.12, and 6.13) complement each other by showing that rwPHP(PLS) is the tightest complexity class in all these results.

Recall that an rwPHP(PLS) instance consists of  $(f, \{I_y\}_{y \in [2M]}, \{g_y\}_{y \in [2M]})$ , where:

- $f : [M] \rightarrow [2M]$  is a purported “surjection”;
- for each  $y \in [2M]$ ,  $I_y := (L, S_y)$  is an instance of ITER, where  $S_y : [L] \rightarrow [L]$ ; and
- $g_y : [L] \rightarrow [M]$  maps solutions of  $I_y$  to integers in  $[M]$ .

We now state and prove the main theorem of this subsection.

**Theorem 5.4.** *There is a universal constant  $C \geq 2$  such that the following holds. Let  $L, M \geq 1$  be the parameters of rwPHP(PLS) instances and  $n \geq 1$ .*

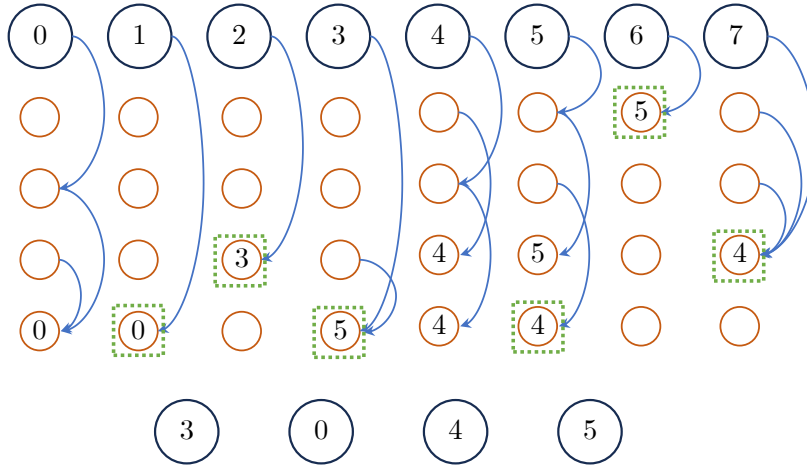
*For every unsatisfiable CNF formula  $F$  over  $n$  variables and parameter  $s_F \geq C \cdot (nLM + |F|)$  such that every resolution refutation of  $F$  requires more than  $s_F$  clauses, there is a decision tree reduction of block-depth  $O(n)$  from a rwPHP(PLS) instance to a  $\text{REFUTER}(s(F \vdash_{\text{Res}} \perp) \leq s_F)$  instance.*

*Proof.* Let  $(f, \{I_y\}_{y \in [2M]}, \{g_y\}_{y \in [2M]})$  be an instance of rwPHP(PLS) and we will reduce it to an instance of  $\text{REFUTER}(s(F \vdash_{\text{Res}} \perp) \leq s_F)$ . Our goal is to construct a size- $s_F$  resolution refutation  $\Pi$  for  $F$  such that any illegal derivation in  $\Pi$  corresponds to a valid solution to the rwPHP(PLS) instance.

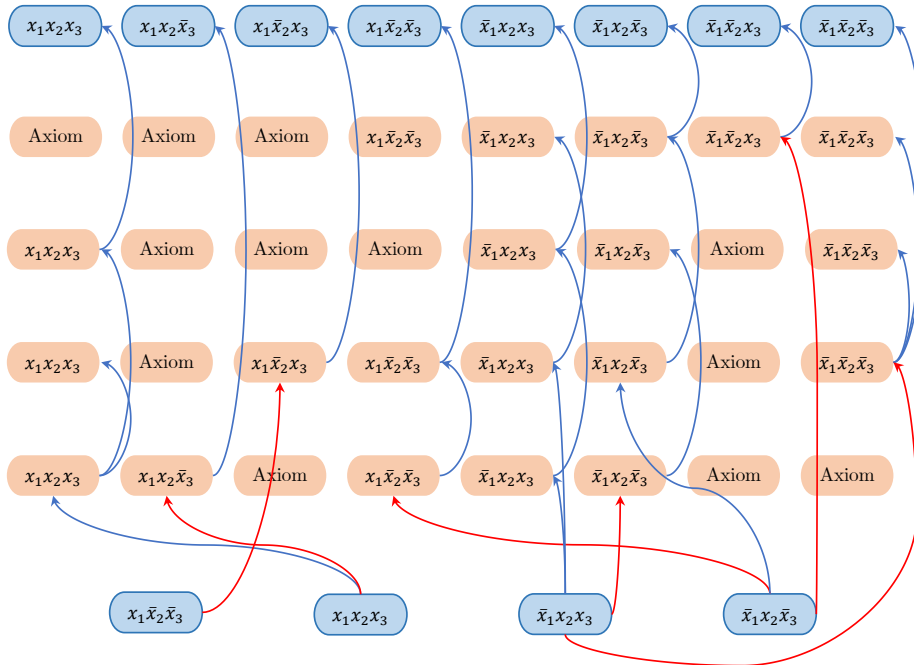
The nodes in  $\Pi$  are partitioned into  $n + 1$  layers, numbered from layer 0 to layer  $n$ . Each layer  $t \in [n + 1]$  has either two rows or one row: When layer  $t$  has two rows, we denote the nodes in the first row by  $\{D_{(y,a)}^t\}$  and those in the second row by  $\{E_i^t\}$ ; when layer  $t$  has only one row, we denote the nodes by  $\{E_i^t\}$ . (Therefore,  $\{E_i^t\}$  always denote the *last* row of layer  $t$ .) After all these  $n + 1$  layers of clauses, we put the axioms of  $F$  at the very end. It is easy to translate a resolution refutation in this layout into one in the format of Definition 3.5 by decision trees of block-depth 1.

**The construction.** The layer 0 has one node  $E_0^0 := \perp$ . For each  $t$  from  $1, \dots, n$ :

1. Let  $E_0^{t-1}, \dots, E_{k-1}^{t-1}$  be the nodes on the last row of layer  $t - 1$ ; we will always guarantee that  $k \leq M$ .
2. **Case 1:**  $2k \leq M$ . In this case, layer  $t$  will only have one row of nodes, defined as follows. For every node  $E_i^{t-1}$  on layer  $t - 1$ , we generate 2 nodes  $E_{2i}^t$  and  $E_{2i+1}^t$  on layer  $t$ , where the clauses written are  $E_{2i}^t = E_i^{t-1} \vee x_t$  and  $E_{2i+1}^t = E_i^{t-1} \vee \bar{x}_t$ . We define  $E_i^{t-1}$  to be resolved from  $E_{2i}^t$  and  $E_{2i+1}^t$ .
3. **Case 2:**  $2k > M$ .
  - (a) First, prepare  $2M$  nodes  $D_{(0,0)}^t, D_{(1,0)}^t, \dots, D_{(2M-1,0)}^t$ . It would be instructive to think of  $\{D_{(y,a)}^t : a \in [L]\}$  for each fixed  $y$  as a chain and we are now preparing the heads of these  $2M$  chains. In what follows, we denote  $C_i = D_{(i,0)}^t$  for ease of notation.  
For each  $i \in [k]$ , let  $C_{2i} = E_i^{t-1} \vee x_t$ ,  $C_{2i+1} = E_i^{t-1} \vee \bar{x}_t$ , and define  $E_i^{t-1}$  to be resolved from  $C_{2i}$  and  $C_{2i+1}$ . We make sure that there are *exactly*  $2M$  clauses on the first row by making several copies of  $C_{2k-1}$ : for each  $i \in \{2k, \dots, 2M - 1\}$ , let  $C_i = C_{2k-1}$ .
  - (b) Generate  $M$  nodes on the second row of layer  $t$ : for every  $i \in [M]$ , let  $E_i^t := C_{f(i)}$ . (Intuitively, if  $f : [M] \rightarrow [2M]$  were a surjection, then every node in  $\{C_i\}_{i \in [2M]}$  would appear in  $\{E_i^t\}_{i \in [M]}$ .)



(a) This is an rwPHP(PLS) instance with  $M = 4$ , compressing the top  $2M$  elements to the bottom  $M$  elements. The bottom four points represent the function  $f : [M] \rightarrow [2M]$ ; i.e., in this example,  $f(0) = 3, f(1) = 0, f(2) = 4, f(3) = 5$ . Every column is a PLS instance (every vertex without an outgoing edge has a self-loop). Every sink is a solution of the corresponding PLS instance, on which we have  $g_y : [L] \mapsto [M]$ . The number on every sink represents  $f(g_y(\cdot))$ , which if different from  $y$ , would be a solution of the whole rwPHP(PLS) instance. In this figure, every solution is marked with a dotted green box.



(b) Part of the constructed resolution derivation II. Initially, we have  $2M = 8$  clauses. At the bottom, we have  $M = 4$  clauses, which exactly correspond to the 3rd, 0th, 4th, and 5th clauses above. For every node  $a$  in a PLS instance, if it is a self-loop, then we let the clause be a weakening from some axiom (and it would never be used again). If  $a$  is a solution of the PLS instance, we let it be the weakening of clause  $g_y(a)$  at the bottom. Otherwise, we let it be the weakening of  $S_y(a)$ . All blue arrows are valid weakenings and *all red arrows are invalid weakenings*. The invalid weakenings here will be the (only) solutions to the refuter problem.

Figure 4: The gadget to embed an rwPHP(PLS) instance.

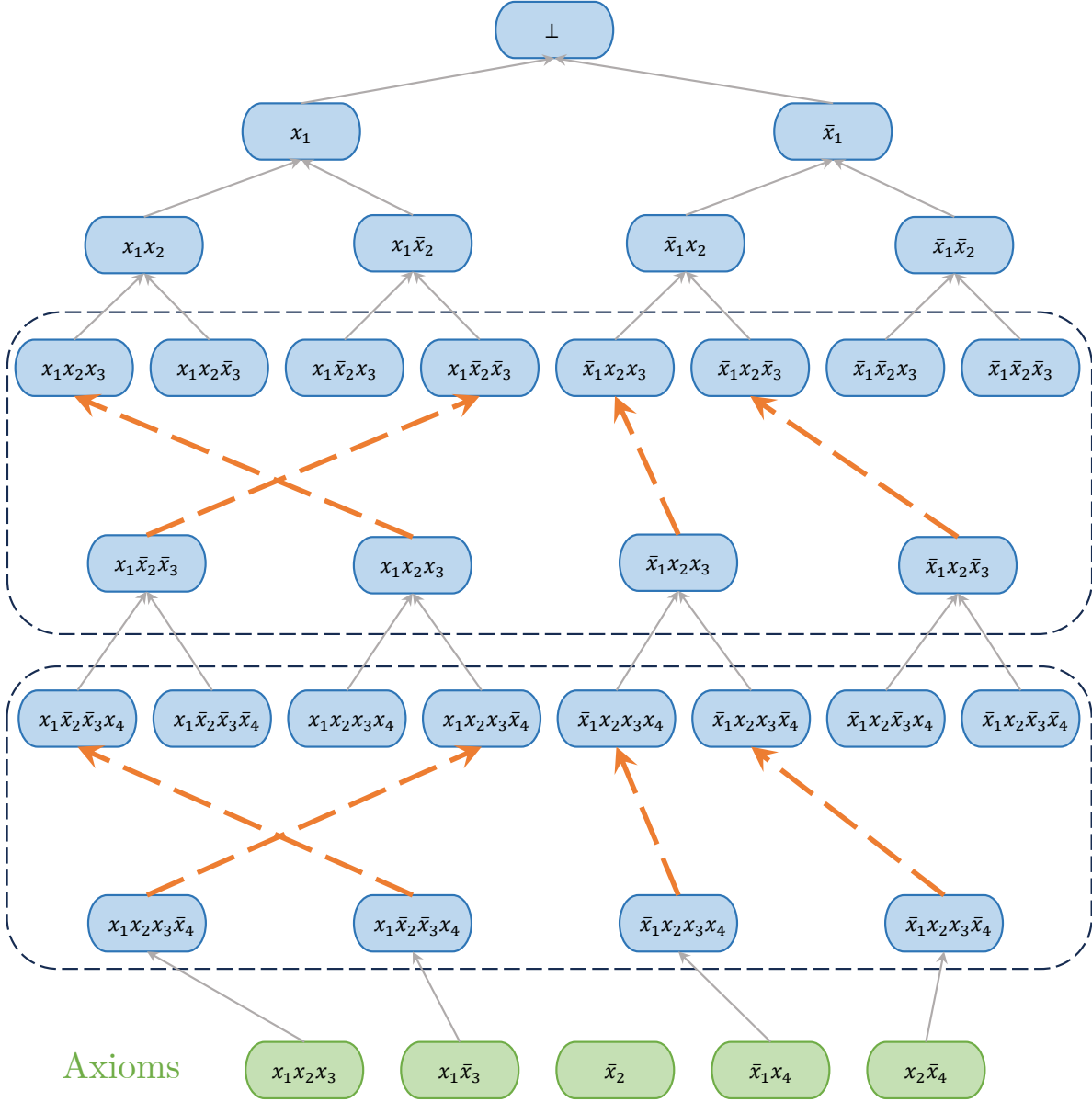


Figure 5: An illustration of our reduction from rwPHP(PLS) to the size refuter problem. All gray arrows are valid resolution derivations (and the last layer is weakening from axioms). Every dashed box uses the gadget to embed an rwPHP(PLS) instance that enforces every layer to have at most  $2M$  clauses. Thus the only possible invalid derivations are those inside the gadget which, once found, would directly imply a solution of the original rwPHP(PLS) instance. The overall reduction will produce a purported resolution refutation of size  $O(nLM + |F|)$ .

- (c) Now, for each  $y \in [2M]$ , we “link” the node  $C_y = D_{(y,0)}^t$  to their corresponding  $E_{f^{-1}(y)}^t$  on the second row, using the ITER instance  $I_y$ . Recall that for each  $y \in [2M]$ ,  $I_y$  consists of a function  $S_y : [L] \rightarrow [L]$ , and solutions of  $I_y$  are those  $a \in [L]$  such that

$$\text{either } S_y(a) < a \text{ or } (S_y(a) > a \text{ and } S_y(S_y(a)) = S_y(a)).$$

As a special case, if  $a = 0$  and  $S_y(0) = 0$ , then  $0$  also counts as a solution.

Every clause on the chain  $\{D_{(y,a)}^t : a \in [L]\}$  will be equal to  $C_y = D_{(y,0)}^t$  except those on a

“junk” node  $a$  such that  $S_y(a) = a$  (see [Case 2 \(c\) ii.](#)); these clauses will be weakenings of each other, and the instance  $I_y$  dictates the structure of the weakening relationship. For every  $a \in [L]$ , the node  $D_{(y,a)}^t$  is defined as follows:

- i. If  $a$  is a solution of  $I_y$ , then  $g_y(a)$  is a purported pre-image of  $y$ . The clause written on  $D_{(y,a)}^t$  is equal to  $C_y$  and we define it to be a weakening of  $E_{g_y(a)}^t$ . (Note that  $E_{g_y(a)}^t = C_{f(g_y(a))}$  by definition, meaning that if the weakening from  $E_{g_y(a)}^t$  to  $D_{(y,a)}^t$  is an illegal derivation, then  $f(g_y(a)) \neq y$ .)
- ii. If  $a$  is not a solution and  $S_y(a) = a$ , then the clause written on  $D_{(y,a)}^t$  is defined to be the weakening of an arbitrary axiom in  $F$  (say the first axiom). The node  $D_{(y,a)}^t$  is considered a “junk” node and will never be used later.
- iii. Otherwise, we have  $S_y(a) > a$ . The clause written on  $D_{(y,a)}^t$  is equal to  $C_y$  and we define it to be a weakening of  $D_{(y,S_y(a))}^t$ .

After constructing all these nodes above, we put the axioms of  $F$  at the very end. Each clause  $E_i^n$  in the last row of layer  $n$  will be a weakening of some axiom in  $F$ . In particular, note that each  $E_i^n$  has exactly  $n$  literals (this can be easily seen from induction) and therefore is satisfied by exactly one assignment  $\alpha_i$ . Recall that  $F$  is an unsatisfiable CNF formula, so for each  $E_i^n$ , there exists an axiom  $A$  in  $F$  such that  $\alpha_i$  falsifies  $A$ ; hence we can define  $E_i^n$  to be a weakening of  $A$ .

This finishes the construction.

The above construction gives a resolution refutation  $\Pi$  for  $F$  that has size  $s_F := O(nLM + |F|)$ . The only place in  $\Pi$  where illegal derivations might occur is in [Case 2 \(c\) i.](#) when we define  $D_{(y,a)}^t$  to be a weakening of  $E_{g_y(a)}^t$ . If this is an illegal derivation, then  $f(g_y(a)) \neq y$ , which means that we have found a valid solution for the rwPHP(PLS) instance. Therefore, the above construction is a correct reduction from rwPHP(PLS) to  $\text{REFUTER}(s(F \vdash_{\text{Res}} \perp) > s_F)$ , as long as  $s_F > C \cdot (nLM + |F|)$  for some large universal constant  $C$ .

Finally, we analyze the query complexity of this reduction. It suffices to show that every node  $D_{(y,a)}^t$  and  $E_i^t$  can be computed in block-depth  $O(n)$  from the input rwPHP(PLS) instance. Note that to compute one node, we need to calculate both its origin (i.e., resolved or weakening from which node) and the clause written on it. We use induction on  $t$  to show that every clause in layer  $t$  can be computed in block-depth  $c \cdot (t + 1)$  for some universal constant  $c$ . Fix a layer  $t$  and we argue as follows.

- **(Base case)** If layer  $t$  contains only one row, then we can read off the clause  $E_i^t$  from the binary representation of  $i$ ; the node  $E_i^t$  is always resolved from  $E_{2i}^{t+1}$  and  $E_{2i+1}^{t+1}$  (if layer  $t + 1$  also contains only one row) or  $C_{2i}^{t+1}$  and  $C_{2i+1}^{t+1}$  (otherwise).
- **(Induction step)** If layer  $t$  contains two rows, then we argue as follows.
  1. For  $i < 2k$ , depending on the parity of  $i$ , we have that the clause written on  $D_{(i,0)}^t$  is either  $E_{\lfloor i/2 \rfloor}^{t-1} \vee x_t$  or  $E_{\lfloor i/2 \rfloor}^{t-1} \vee \bar{x}_t$ . For  $i \geq 2k$ , the clause written on  $D_{(i,0)}^t$  is always equal to  $D_{(2k-1,0)}^t$ . For every  $i \in [2M]$  and  $a \in [L]$ , the clause written on  $D_{(i,a)}^t$  is either equal to the clause written on  $D_{(i,0)}^t$ , or equal to some axiom of  $F$ , and this can be decided in block-depth 2 (see [Case 2 \(c\) ii.](#)). Since it takes block-depth  $ct$  to compute  $E_{\lfloor i/2 \rfloor}^{t-1}$ , it takes block-depth  $ct + 2$  to compute the clause written on  $D_{(i,a)}^t$ .
  2. Every  $D_{(y,a)}^t$  (for  $y \in [2M]$  and  $a \in [L]$ ) belongs to one of the following three cases:
    - if  $a$  is a solution of  $I_y$ , then  $D_{(y,a)}^t$  is a weakening of  $E_{g_y(a)}^t$ ;
    - if  $a \neq 0$  and  $S_y(a) = a$ , then  $D_{(y,a)}^t$  is a weakening of some axiom in  $F$  and is a “junk” node;

– otherwise,  $D_{(y,a)}^t$  is a weakening of  $D_{(y,S_y(a))}^t$ .

Therefore, we can use  $O(1)$  additional block-depth to determine all information regarding  $D_{(y,a)}^t$ .

3. Let  $i \in [M]$ , then  $E_i^t = D_{(f(i),0)}^t$ , and  $E_i^t$  is either resolved from  $C_{2i}^{t+1}$  and  $C_{2i+1}^{t+1}$  (when  $t < n$ ) or a weakening of some axiom (when  $t = n$ ). This can be computed in constant additional block depth.

It follows that  $\Pi$  can be computed from our input  $\text{rwPHP(PLS)}$  instance in block-depth  $O(n)$ .  $\square$

**Corollary 5.5.**  $\text{REFUTER}(s(\text{PHP}_{(n+1) \rightarrow n} \vdash_{\text{Res}} \perp) \leq 1.01^n)$  is complete for  $\text{rwPHP(PLS)}$ .

*Proof Sketch.* By combining [Theorem 4.16](#) and [Theorem 5.4](#). Note that the reductions have  $\text{poly}(n)$  block-depth and each block contains  $\text{poly}(n)$  bits, therefore they are polynomial-time (many-one) reductions.  $\square$

The above hardness result in  $\text{TFNP}$  can be interpreted as a reversal result in bounded reverse mathematics as well. To state this reversal result, we define the following two families of  $\forall\Sigma_1^b(\alpha)$ -sentences. For  $\text{PV}(\alpha)$  function symbols  $F, I, G$ , let  $\text{rwPHP(PLS)}(F, I, G)$  denote the natural  $\forall\Sigma_1^b(\alpha)$ -sentence expressing the existence of a solution for the  $\text{rwPHP(PLS)}$ -instance defined by  $(F, I, G)$ :

- For every auxiliary input  $z$  and every  $t, L$ , there exists  $y \in [2t]$  and  $\text{ans} \in [L]$  such that  $\text{ans}$  is a PLS solution for the  $\text{ITER}$  instance  $I_{z,y} : [L] \rightarrow [L]$  and that  $(G_{z,y}(\text{ans}) > t$  or  $F_z(G_{z,y}(\text{ans})) \neq y$ ).

Similarly, let  $\text{LB}_{\text{PHP}}^{\text{Res}}(\text{PV}(\alpha))$  denote the family of  $\forall\Sigma_1^b(\alpha)$ -sentences consisting of

$$\forall n \in \text{Log} \forall z \exists i \in [1.01^n] \text{mistake}_{\text{PHP}}(n, M, z, i)$$

for every  $\text{PV}(\alpha)$  function symbol  $M(i, z)$  (here  $z$  is a parameter).

**Theorem 5.6.** For every  $\text{PV}(\alpha)$  function symbols  $F, I, G$ ,

$$\text{PV}(\alpha) + \text{LB}_{\text{PHP}}^{\text{Res}}(\text{PV}(\alpha)) \vdash \text{rwPHP(PLS)}(F, I, G).$$

*Proof Sketch.* Argue in  $\text{PV}(\alpha)$ . Let  $\Pi$  be the purported resolution proof for  $\text{PHP}$  as constructed in the proof of [Theorem 5.4](#) from  $(F, I, G)$ , then  $\Pi$  can be expressed as a  $\text{PV}(\alpha)$  function symbol (that depends on  $F, I$ , and  $G$ ). From  $\text{LB}_{\text{PHP}}^{\text{Res}}(\text{PV}(\alpha))$ , we know that there exists an illegal derivation in  $\Pi$ . This illegal derivation can only occur in [Case 2 \(c\) i.](#), and hence it points to a weakening from some  $D_{(y,\text{ans})}^t$  to some  $E_{g_y(\text{ans})}^t$ . This means the existence of a solution  $(y, \text{ans})$  of the  $\text{rwPHP(PLS)}$ -instance  $(F, I, G)$ .  $\square$

We remark that like [Theorem 5.4](#), the proof of the above theorem does not depend on the hard tautology being  $\text{PHP}$ .

We finish this section by the following nice-looking characterization of  $\forall\Sigma_1^b$ -consequences (i.e., provably total NP search problems) of  $\mathbb{T}_2^1 + \text{dwPHP(PV)}$ :

**Corollary 5.7.** 1.  $\text{REFUTER}(s(\text{PHP}_{(n+1) \rightarrow n} \vdash_{\text{Res}} \perp) \leq 1.01^n)$  is complete for the class of NP search problems provably total in  $\mathbb{T}_2^1 + \text{dwPHP(PV)}$ .

2. A  $\forall\Sigma_1^b(\alpha)$ -sentence is provable in the theory  $\mathbb{T}_2^1(\alpha) + \text{dwPHP(PV}(\alpha))$  if and only if it is provable in the theory  $\text{PV}(\alpha) + \text{LB}_{\text{PHP}}^{\text{Res}}(\text{PV}(\alpha))$ .

## 6 Refuters for Other Formulas

This section presents additional upper bounds for the refuter problems associated with resolution lower bounds. We start with a *universal* PLS upper bound for width refuters, showing that any resolution width lower bound *that is true* can be refuted in non-uniform PLS. Then, we provide further examples of resolution lower bounds proven by “random restriction + width lower bounds” and show that the refuter problems for these lower bounds are in rwPHP(PLS). In particular, we present the following three classic resolution lower bounds and show that the refuter problems for all of them are in rwPHP(PLS):

- (a) size-width tradeoffs from XOR-lifting [DR03, Kra11a] (Section 6.2);
- (b) exponential size lower bounds for the Tseitin formulas [Urq87, Sch97] (Section 6.3); and
- (c) exponential size lower bounds for random  $k$ -CNFs [CS88, BP96] (Section 6.4).

We believe that the case of random  $k$ -CNFs is especially compelling: the *vast majority* of resolution lower bounds have refuters in rwPHP(PLS)!

### 6.1 Universal Refuters for *Every* Narrow Resolution Proof

This subsection shows a very general result: For *every* (possibly non-uniform) family of unsatisfiable CNFs  $\mathcal{F} = \{F_n\}$  and *every* sequence of integers  $\{w_n\}$ , if for every  $n \in \mathbb{N}$ ,  $w_n$  is indeed a resolution width lower bound for  $F_n$ , then the refuter problem corresponding to this width lower bound is in PLS under non-uniform decision tree reductions.

We note that such a membership result is *inherently* non-uniform since it is crucial to consider algorithms with unlimited *computational* power. For example, in general, it is not obvious how to decide if  $w_n$  is a valid resolution width lower bound for  $F_n$  (although it is certainly computable with unlimited computational power). In fact, even checking if  $F_n$  is unsatisfiable is itself NP-complete. On the other hand, even though these two tasks are computationally hard, they only require querying at most  $\text{poly}(n)$  bits of the given resolution proof. Thus, we can still consider these refuter problems in  $\text{TFNP}^{\text{dt}}$  and study its query complexity in the non-uniform setting.

**Theorem 6.1.** *Let  $\mathcal{F}$  be any (possibly non-uniform) family of unsatisfiable CNFs with polynomially many clauses. Let  $w_0 = w(\mathcal{F} \vdash_{\text{Res}} \perp)$ . Then there exists a (non-uniform) decision-tree reduction of block-depth 2 from  $\text{REFUTER}(w(\mathcal{F} \vdash_{\text{Res}} \perp) < w_0)$  to ITER.*

*Proof.* Consider any instance of  $\text{REFUTER}(w(\mathcal{F} \vdash_{\text{Res}} \perp) < w_0)$ . Recall from Definition 3.5 that the instance is a purported resolution refutation  $\Pi$  that consists of clauses  $C_{-k}, \dots, C_{-1}, C_0, \dots, C_{L-1}$  where  $C_{-k}, \dots, C_{-1}$  are the axioms of  $\mathcal{F}$  and  $C_{L-1} = \perp$ . Also, recall that we syntactically ensure the width of  $\Pi$  is  $< w_0$  by only allocating  $w_0 - 1$  literals for each clause. The key point in the reduction is that, for any clause  $C_i$  that is resolved from  $C_{j_1}$  and  $C_{j_2}$ , if  $\text{width}(F \vdash_{\text{Res}} C_i) \geq w_0$ , then either  $\text{width}(F \vdash_{\text{Res}} C_{j_1}) \geq w_0$ , or  $\text{width}(F \vdash_{\text{Res}} C_{j_2}) \geq w_0$ .

The length of the reduced reversed ITER instance is exactly  $L$ . Next, we define the successor pointers: for every  $i \in [L]$ , let  $C_i$  be the  $i$ -th clause and  $C_{j_1}$  and  $C_{j_2}$  with  $j_1 < j_2 < i$  be the two clauses from which  $C_i$  is resolved, then

$$S(i) := \begin{cases} i & \text{if } \text{width}(F \vdash_{\text{Res}} C_i) < w_0; \\ j_1 & \text{if } \text{width}(F \vdash_{\text{Res}} C_{j_1}) \geq w_0; \\ j_2 & \text{otherwise.} \end{cases}$$

Clearly, this is a query-efficient reduction with block-depth 2. It is not time-efficient because it needs to compute whether  $\text{width}(F \vdash_{\text{Res}} C) < w_0$  for some clauses  $C$ .

To show correctness, we consider any possible solution of the constructed reversed ITER. For any  $i$  such that  $S(i) > i$ , we have either  $j_1 > i$  or  $j_2 > i$ , which means that  $C_i$  is an invalid derivation. Now

consider any  $i$  such that  $S(i) < i$  and  $S(S(i)) = S(i)$ . Since  $S(i) < i$ , we have that  $\text{width}(F \vdash_{\text{Res}} C_i) \geq w_0$ . Since  $S(S(i)) = S(i)$ , we have both  $\text{width}(F \vdash_{\text{Res}} C_{j_1}) < w_0$  and  $\text{width}(F \vdash_{\text{Res}} C_{j_2}) < w_0$ . Thus, the resolution step from  $C_{j_1}$  and  $C_{j_2}$  to  $C_i$  must be an invalid derivation. This finishes the proof.  $\square$

Note that [Theorem 5.1](#) already shows a universal PLS-hardness, which even holds for uniform reduction. Combining the the PLS-membership ([Theorem 6.1](#)) above, we have the following corollary.

**Corollary 6.2.** *Let  $\mathcal{F}$  be any (possibly non-uniform) family of unsatisfiable CNFs with polynomially many clauses. Let  $w_0 := w(\mathcal{F} \vdash_{\text{Res}} \perp)$ . Then  $\text{REFUTER}(w(\mathcal{F} \vdash_{\text{Res}} \perp) < w_0)$  is PLS-complete under (non-uniform) decision tree reductions.*

## 6.2 Refuters for XOR-Lifted Lower Bounds

We show that for a large family of resolution lower bounds proved by *lifting theorems*, their corresponding refuter problems are in  $\text{rwPHP}(\text{PLS})$ .

Given an unsatisfiable CNF  $F$  which is hard for a “weak” proof system, a *lifting theorem* produces another unsatisfiable CNF  $F'$  (typically by composing  $F$  with some *gadgets*) that is hard for a “stronger” proof system. Lifting is a very influential technique for proving lower bounds in proof complexity, see e.g. [\[HN12, GP18b, dRNV16, GGKS20, dRMN+20\]](#). This subsection examines one of the simplest lifting theorems for proving lower bounds for resolution, which originated from the technique of “relativization” [\[DR03, Kra11a\]](#) (see also [\[Kra19, Section 13.2\]](#)).

Let  $F(z_1, z_2, \dots, z_n)$  be an unsatisfiable CNF. Roughly speaking, the CNF  $F \circ \text{XOR}$  is obtained by replacing each variable  $z_i$  with  $x_i \oplus y_i$ , where  $x_i$  and  $y_i$  are new variables corresponding to  $z_i$ . More formally, the formula  $F \circ \text{XOR}$  takes  $2n$  Boolean variables  $x_1, x_2, \dots, x_n$  and  $y_1, y_2, \dots, y_n$  as inputs. Denoting  $z_i^b = z_i$  if  $b = 1$  and  $\bar{z}_i$  if  $b = 0$ ; each width- $d$  clause

$$z_{i_1}^{b_1} \vee z_{i_2}^{b_2} \vee \dots \vee z_{i_d}^{b_d}$$

becomes a set of  $2^d$  width- $2d$  clauses

$$\left\{ \left( x_{i_1}^{r_1 \oplus 1} \vee y_{i_1}^{b_1 \oplus r_1} \right) \vee \left( x_{i_2}^{r_2 \oplus 1} \vee y_{i_2}^{b_2 \oplus r_2} \right) \vee \dots \vee \left( x_{i_d}^{r_d \oplus 1} \vee y_{i_d}^{b_d \oplus r_d} \right) \right\}_{r_1, r_2, \dots, r_d \in \{0,1\}}.$$

A classical lifting theorem states that if  $F$  requires large resolution *width*, then  $F \circ \text{XOR}$  requires large resolution *size*. Here, the “weak” proof system is *narrow* resolution and the “strong” proof system is *short* resolution. More formally:

**Theorem 6.3.** *Let  $F$  be an unsatisfiable CNF that requires resolution width  $\geq w$ , then  $F \circ \text{XOR}$  requires resolution size  $\geq 2^{w/3}$ .*

The classical proof of this theorem goes through a random restriction argument. Let  $\Pi$  be a purported resolution proof of  $F \circ \text{XOR}$  of length  $L < 2^{w/3}$ . Consider a random restriction  $\rho$  as follows: For each index  $i$ , with probability  $1/2$ , we set  $\rho_{x_i} = 0/1$  uniformly at random and  $\rho_{y_i} = *$ ; otherwise, we set  $\rho_{y_i} = 0/1$  uniformly at random and  $\rho_{x_i} = *$ . By the construction above,  $\Pi|_\rho$  is a resolution proof of  $F$  up to substituting some variables by their negations, for any  $\rho$  in the support. Moreover, for any clause  $C \in \Pi$  of width at least  $t$ ,  $C$  is killed by a random restriction  $\rho$  (i.e.,  $C|_\rho \equiv 1$ ) w.p. at least  $1 - 2^{-\Omega(t)}$ . By a union bound over all  $L < 2^{w/3}$  clauses in  $\Pi$ , it follows that there is a random restriction  $\rho$  killing every clause of width  $> w$  in  $\Pi$ . Therefore,  $\Pi|_\rho$  is a resolution refutation for  $F$ , contradicting the width lower bound for  $F$ .

To obtain a reduction to  $\text{rwPHP}(\text{PLS})$ , it would be helpful to rephrase the above proof as a *compression* argument:

*Proof of Theorem 6.3.* Let  $\mathcal{R}$  be the space of the random restrictions in the above proof. Each  $\rho \in \mathcal{R}$  can be described in  $2n$  bits:

- For each index  $i$ , if  $\rho_{x_i} = *$ , then we write down  $0y_i$  (the first bit being 0 indicates that  $x_i$  is set to  $*$ , and the second bit encodes  $y_i$ ); otherwise we write down  $1x_i$ .

We call the above encoding the *standard encoding* of a restriction; this encoding is a bijection between  $\mathcal{R}$  and  $\{0, 1\}^{2n}$ , showing that  $|\mathcal{R}| = 4^n$ .

In contrast, if  $C$  is a clause of width  $w$  and  $\rho \in \mathcal{R}$  is a restriction such that  $C|_\rho \neq 1$ , then given  $C$ , such a  $\rho$  can be described in  $(\log_2 3)w + 2(n - w) < 2n$  bits. This is because for each literal in  $C$  (say  $x_i, y_i, \bar{x}_i$ , or  $\bar{y}_i$ ), if this literal is not simplified to 1, then there are only 3 possible choices for  $(\rho_{x_i}, \rho_{y_i})$ ; for example, if this literal is  $x_i$ , then  $(\rho_{x_i}, \rho_{y_i})$  might be one of  $(0, *)$ ,  $(*, 0)$ , or  $(*, 1)$ , but never  $(1, *)$ . We call this the *short encoding* of  $\rho$  w.r.t.  $C$ ; note that this encoding only works when  $C|_\rho \neq 1$ .

Now, let  $\Pi = (C_0, C_1, \dots, C_{L-1})$  be a resolution refutation of  $F \circ \text{XOR}$  with  $L < 2^{w/3}$  clauses. Let  $f : [L] \times [3^w 4^{n-w}] \rightarrow [4^n]$  be the function that on input  $(i, \rho')$ , where  $i \in [L]$  and  $\rho'$  is the short encoding of a restriction w.r.t.  $C_i$ , outputs the standard encoding of  $\rho'$  in  $\{0, 1\}^{2n}$ . Since

$$L \times 3^w 4^{n-w} \leq 2^{w/3} \cdot 4^n (3/4)^w < 0.99 \cdot 4^n \text{ (whenever } w \geq 1),$$

it follows from the *dual weak pigeonhole principle* that there exists a  $\rho \in \{0, 1\}^{2n}$  outside the range of  $f$ . This restriction  $\rho$  simplifies  $\Pi$  into a width- $w$  resolution proof of  $F$ .

In conclusion, if there is a resolution refutation of  $F \circ \text{XOR}$  with  $< 2^{w/3}$  clauses, then *by the dual weak pigeonhole principle*, there is a resolution refutation of  $F$  with width  $< w$ , contradicting the assumed hardness of  $F$ .  $\square$

Now we are ready to show the following result: for every unsatisfiable CNF of the form  $F \circ \text{XOR}$  whose resolution size lower bound can be derived from [Theorem 6.3](#), the refuter problem for this resolution size lower bound is in  $\text{rwPHP}(\mathcal{P})$ , where  $\mathcal{P}$  corresponds to the refuter problem for the width lower bound for  $F$ . Since the refuter problem corresponding to *every* resolution width lower bound admits a non-uniform reduction to PLS ([Corollary 6.2](#)), the refuter problems corresponding to size lower bounds for  $F \circ \text{XOR}$  non-uniformly reduce to  $\text{rwPHP}(\text{PLS})$  as well. Even if we restrict ourselves to uniform reductions, the refuter problems for many interesting width lower bounds reduce to PLS (such as [Theorem 4.3](#)), thus the refuter problems for size lower bounds for the corresponding lifted CNFs also reduce to  $\text{rwPHP}(\text{PLS})$ .

**Theorem 6.4.** *Let  $\{F_n\}$  be a family of unsatisfiable CNFs,  $w(n)$  be a width lower bound for  $F_n$ , and  $\mathcal{P}$  denote the problem  $\text{REFUTER}(w(F_n) > w(n))$ . Then there is a decision tree reduction from  $\text{REFUTER}(s(F_n \circ \text{XOR}) < 2^{w(n)/3})$  to  $\text{rwPHP}(\mathcal{P})$  with block-depth 1.*

*Proof.* Let  $\Pi$  be the input instance of  $\text{REFUTER}(s(F_n \circ \text{XOR}) < 2^{w(n)/3})$ . That is,  $\Pi = (C_0, C_1, \dots, C_{L-1})$  is a purported resolution refutation of  $F_n \circ \text{XOR}$  with  $L < 2^{w(n)/3}$  clauses, and we want to find an invalid derivation in  $\Pi$ .

Let  $f : [0.99N] \rightarrow [N]$  be the function defined in the proof of [Theorem 6.3](#), where  $N := 4^n$ . That is, given a pair  $(i, \rho')$  where  $i \leq L$  and  $\rho'$  is the short encoding of a restriction w.r.t.  $C_i$ ,  $f(i, \rho')$  is the standard encoding of this restriction. The range of  $f$  consists of (the standard encodings of) all *bad* restrictions, i.e., those that do *not* simplify  $\Pi$  to a width- $w$  resolution refutation.

Given any restriction  $\rho \in \{0, 1\}^{2n}$ , let  $\Pi|_\rho$  denote the restriction of  $\Pi$  under  $\rho$  where we force every clause to have width at most  $w$ ;  $\Pi|_\rho$  is a purported width- $w$  resolution refutation for  $F_n$ . In particular, for each  $(\rho, i)$ , the  $i$ -th clause of  $\Pi|_\rho$  is equal to the restriction of  $C_i$  under  $\rho$ , truncated at width  $w$ . Note that if  $F_n$  indeed requires resolution width  $> w$ , then  $\Pi|_\rho$  must be an *invalid* resolution refutation of  $F_n$ . Suppose that the  $i$ -th clause in  $\Pi|_\rho$  is derived illegally, then it could be for the following two reasons:

- Either the derivation of  $C_i$  in  $\Pi$  is already illegal;
- or the width of  $C_i|_\rho$  is actually  $> w$  and the  $i$ -th clause in  $\Pi|_\rho$  is illegal because it was truncated.

Let  $g'_{\rho,i}$  denote the short encoding of  $\rho$  w.r.t.  $C_i$ , and define  $g_{\rho,i} := (i, g'_{\rho,i})$ . In the second case, we have a clause  $C_i$  and a restriction  $\rho$  such that  $C_i|_{\rho} \neq 1$  (in fact, the width of  $C_i|_{\rho}$  is large), thus the short encoding makes sense and, indeed,  $f(g_{\rho,i}) = \rho$ . In any case, if the short encoding does not make sense (i.e.,  $C_i|_{\rho} = 1$ ), we can set  $g_{\rho,i}$  arbitrarily.

Now we have all the ingredients needed in our reduction from the problem of finding an invalid derivation in  $\Pi$  to  $\text{rwPHP}(\mathcal{P})$ :

1. a purported ‘‘surjection’’  $f : [0.99N] \rightarrow [N]$ ;
2. a  $\mathcal{P}$  instance  $\Pi|_{\rho}$  for every  $\rho \in [N]$ ;
3. for every  $\rho \in [N]$  and every solution  $i$  of  $\Pi|_{\rho}$  (as a  $\mathcal{P}$  instance), a number  $g_{\rho,i}$  pointing to a purported pre-image in  $f^{-1}(\rho)$ .

Every entry  $f(i, \rho')$ ,  $\Pi|_{\rho}(i)$ , and  $g_{\rho,i}$  only depend on  $C_i$  and  $\rho$ , thus are computable by a decision tree of block-depth 1.

A solution of the above  $\text{rwPHP}(\mathcal{P})$  instance consists of a restriction  $\rho$  and a solution  $i$  of  $\Pi|_{\rho}$  such that  $f(g_{\rho,i}) \neq \rho$ . In this case, the derivation of  $C_i$  in  $\Pi$  must be invalid. That is, given a solution of the  $\text{rwPHP}(\mathcal{P})$  instance, we can find an invalid derivation of  $\Pi$  by a decision tree of block-depth 1.  $\square$

### 6.3 Refuters for Tseitin Formulas

**Tseitin formulas.** Let  $G = (V, E)$  be a undirected connected graph, where each vertex  $v \in V$  is associated with a value  $\tau(v) \in \{0, 1\}$ , and each edge  $e \in E$  is associated with a Boolean variable  $x_e$ . The goal is to assign values to each  $x_e$  so that for each vertex  $v \in V$ , the XOR of edge labels incident to  $v$  is equal to  $\tau(v)$ ; that is,

$$\bigoplus_{e \sim v} x_e = \tau(v), \tag{6}$$

where  $e \sim v$  denotes that the edge  $e$  is incident to the vertex  $v$ .

We say  $\tau$  is an *odd-weighted* function if  $\bigoplus_{v \in V} \tau(v) = 1$ . It is not hard to see that the above task is impossible if and only if  $\tau$  is odd-weighted ([Urq87, Lemma 4.1]).

**Definition 6.5.** The *Tseitin formula*  $\text{Tseitin}(G, \tau)$  [Tse83] consists of Equation 6 for every vertex  $v$ . When  $G$  is a  $d$ -regular graph (i.e., every vertex  $v$  is incident to exactly  $d$  edges), we can write Equation 6 as a  $d$ -CNF with  $2^{d-1}$  clauses:

$$\bigwedge_{y_1 \oplus y_2 \oplus \dots \oplus y_d \neq \tau(v)} ((x_{e_1} \neq y_1) \vee (x_{e_2} \neq y_2) \vee \dots \vee (x_{e_d} \neq y_d)), \tag{6'}$$

where  $e_1, e_2, \dots, e_d$  are edges incident to  $v$ .

For every odd-weighted function  $\tau$ ,  $\text{Tseitin}(G, \tau)$  is unsatisfiable; when  $G$  is an *expander* graph,  $\text{Tseitin}(G, \tau)$  becomes hard for resolution.

**Definition 6.6.** Let  $G = (V, E)$  be an undirected graph. For  $S, T \subseteq V$ , denote  $E(S, T)$  as the set of edges in  $E$  with one endpoint in  $S$  and the other endpoint in  $T$ . The *expansion* of  $G$  is defined as:

$$e(G) := \min\{|E(S, V \setminus S)| : S \subseteq V, |V|/3 \leq |S| \leq 2|V|/3\}.$$

This gives rise to a family of popular hard tautologies in proof complexity. The first exponential resolution lower bound for Tseitin formulas was proved by Urquhart [Urq87]; the proof was subsequently simplified by [Sch97, BW01]. We restate the theorem from [BW01] below.

**Theorem 6.7** ([BW01, Theorem 4.4]). *For every undirected connected graph  $G$  and odd-weighted function  $\tau : V \rightarrow \{0, 1\}$ , any resolution refutation of  $\text{Tseitin}(G, \tau)$  contains a clause  $C$  with  $w(C) \geq e(G)$ .*

In this paper, we only consider Tseitin formulas on graphs with constant degree  $d = O(1)$ .

**Width Refuters.** Similar to the Pigeonhole Principle, we first study the width refuter for Tseitin formulas.

**Definition 6.8.** Let  $\text{REFUTER}(w(\text{Tseitin} \vdash_{\text{Res}} \perp) < e(G))$  denote the following problem. The input consists of an undirected connected graph  $G = (V, E)$  on  $n$  vertices with degree  $d = O(1)$ , an odd-weighted assignment  $\tau : V \rightarrow \{0, 1\}$ , a parameter  $e \leq |E|$ , and a purported resolution refutation  $\Pi$  of  $\text{Tseitin}(G, \tau)$  with width less than  $e$ . A valid solution is either of the following:

- an index  $i$  such that the  $i$ -th node in  $\Pi$  is an invalid derivation, or
- a vertex set  $S \subseteq V$  such that  $|V|/3 \leq |S| \leq 2|V|/3$  and  $|E(S, V \setminus S)| < e$ .

(Note: in this  $\text{TFNP}^{\text{dt}}$  problem, we think of  $\text{poly}(n)$ -time algorithms as “efficient”, hence an efficient procedure can read the whole graph  $G$ , verify that  $\tau$  is indeed odd-weighted, or count the number of edges between  $S$  and  $V \setminus S$ . When we calculate block-depth, the inputs  $(G, \tau, e)$  are treated as a single block.)

*Remark 5.* This definition is different from most refuter problems considered in this paper, as it is not for a single family of tautology, and it does not even *guarantee* that the tautology is hard! Instead, it asks to find either an invalid derivation in the purported proof or a *certificate* of the tautology being easy (i.e., a sparse cut in the graph).

We argue that this is a natural definition. Let  $\text{pf}_{\text{Tseitin}}^\alpha(G, \tau, e)$  denote the  $\Pi_1^b(\alpha)$ -sentence “ $\alpha$  encodes a width- $e$  proof of  $\text{Tseitin}(G, \tau)$ ” (note that  $\alpha$  is treated as an oracle, i.e., a second-order object, while  $G, \tau$ , and  $e$  are inputs, i.e., first-order objects). That is,

$$\text{pf}_{\text{Tseitin}}^\alpha(G, \tau, e) := \forall i \text{ Correct}^\alpha(G, \tau, e, i),$$

where  $\text{Correct}^\alpha(G, \tau, e, i)$  expresses that the  $i$ -th step of  $\alpha$ , as a width- $(e-1)$  proof of  $\text{Tseitin}(G, \tau)$ , is correct. Similarly, let  $\text{Expander}(G, e)$  denote the  $\Pi_1^b$ -sentence that  $e(G) \geq e$ . That is,

$$\text{Expander}(G, e) := \forall S \subseteq V (|S| \in [(1/3)|V|, (2/3)|V|] \implies |E[S, V \setminus S]| \geq e).$$

The proof in [BW01] actually shows that  $\text{Expander}(G, e) \implies \neg \text{pf}_{\text{Tseitin}}^\alpha(G, \tau, e)$ , which after rearranging is equivalent to:

$$\exists i \neg \text{Correct}^\alpha(G, \tau, e, i) \vee \exists S (|S| \in [(1/3)|V|, (2/3)|V|] \wedge |E[S, V \setminus S]| < e). \quad (7)$$

It is easy to see that **Definition 6.8** is exactly the  $\text{TFNP}^{\text{dt}}$  problem corresponding to **Equation 7**.

**Theorem 6.9.**  $\text{REFUTER}(w(\text{Tseitin} \vdash_{\text{Res}} \perp) < e(G))$  is PLS-complete.

*Proof.* We will show that there is a (uniform) decision tree reduction of block-width 3 from this problem to ITER.

Let  $G = (V, E)$  be an undirected graph with purported expansion parameter  $e$ . Let  $\tau : V \rightarrow \{0, 1\}$  be an odd-weighted function. Let  $\Pi$  be a purported resolution refutation that consists of clauses  $C_{-k}, \dots, C_{-1}, C_0, \dots, C_{L-1}$ . where  $C_{-k}, \dots, C_{-1}$  are axioms of the unsatisfiable CNF associated with  $G$  and  $\tau$ . Note that we can syntactically require that each  $C_i$  has width at most  $e-1$ . Our goal is to find either an invalid derivation in  $\Pi$  or a witness that the expansion of  $G$  is, in fact, less than  $e$ . In particular, the witness is a vertex set  $S \subseteq V$  such that  $|V|/3 \leq |S| \leq 2|V|/3$  and  $|E(S, V \setminus S)| < e$ .

Similar to before, we first introduce a complexity measure for a clause  $C$ . Let  $v \in V$  be a vertex. We say an assignment  $\alpha$  is  $v$ -critical if  $\alpha$  only falsifies the constraint associated with  $v$  and satisfies all other constraints of the given unsatisfiable CNF. The complexity measure, denoted by  $\text{cri}(C)$ , is defined as follows.

$$\text{cri}(C) := |\{v \in V : \exists v\text{-critical assignment } \alpha \text{ such that } C(\alpha) = 0\}|.$$

Note that  $\text{cri}$  has four important properties:

- $\text{cri}(\perp) = n$ ;

- $\text{cri}(C_i) = 1$  for all  $-k \leq i \leq -1$ , namely,  $\text{cri}(C) = 1$  for all axioms  $C$ ;
- $\text{cri}$  is subadditive with respect to resolution derivation, namely, if  $C$  is resolved from  $A$  and  $B$ , then  $\text{cri}(C) \leq \text{cri}(A) + \text{cri}(B)$ ;
- if  $C$  is obtained from a weakening of  $A$ , then  $\text{cri}(C) \leq \text{cri}(A)$ .

We first show that  $\text{cri}(\cdot)$  can be computed in polynomial time. Then we show that any clause  $C_i$  such that  $n/3 \leq \text{cri}(C_i) \leq 2n/3$  will give us a solution. The PLS-membership follows from that the standard  $1/3$ - $2/3$  trick can be implemented via a reduction to reversed ITER.

**Lemma 6.10.** *For any clause  $C$ ,  $\text{cri}(C)$  can be computed in  $\text{poly}(n)$  time.*

*Proof.* Fix any clause  $C$ . We will enumerate  $v \in V$  and check the existence of  $v$ -critical assignments.

Note that the aimed assignment  $\alpha$  needs to satisfy that  $C(\alpha) = 0$ , so all literals in  $C$  are fixed. For  $\alpha$  being a  $v$ -critical assignment, the constraint associated with  $v$  needs to be falsified. We enumerate an axiom in the constraint associated with  $v$ . Since  $d$  is a constant, there are only  $2^{d-1} = O(1)$  axioms that we need to enumerate.

Fix such an axiom, and set all literals in this axiom to be 0 as well (if setting them to be 0 is not consistent with  $C(\alpha) = 0$ , then skip this axiom and try the next one). Now we have fixed some variables and left other variables free. Let  $\rho \in \{0, 1, *\}^m$  be this partial assignment, where  $m = |E|$ . Note that  $C(\rho) = 0$  and the constraint associated with  $v$  has also been falsified. So we only need to check if there is a complement  $\alpha$  of  $\rho$  such that all other constraints can be satisfied by  $\alpha$ . This reduces to checking whether a system of linear equations over  $\mathbb{F}_2$  has a solution, which can be done in polynomial time.  $\diamond$

Then we show that finding a clause  $C_i$  such that  $\text{cri}(C_i) \in [n/3, 2n/3]$  can be reduced to ITER.

**Reduction to ITER:** The instance of a reversed ITER is defined by the following function  $S : [L] \rightarrow [L]$ . For every  $i \in [L]$ :

- if  $\text{cri}(C_i) < \frac{2n}{3}$ , then  $S(i) = i$ ;
- otherwise, if  $C_i$  is a weakening from  $C_j$ , then let  $S(i) = j$ ;
- Finally, let  $C_i$  be resolved from  $C_j$  and  $C_k$ : If  $\text{cri}(C_j) \geq \text{cri}(C_k)$ , then  $S(i) = j$ ; otherwise  $S(i) = k$ .

It is easy to see that this reduction can be implemented in block-depth 3: for example, if  $C_i$  is resolved from  $C_j$  and  $C_k$ , then one only needs to read the  $i$ -th,  $j$ -th, and  $k$ -th node in the resolution refutation.

Note that when we find any solution  $i$  of this reversed ITER instance, it satisfies  $S(i) < i$  and  $S(S(i)) = i$ . This means  $\text{cri}(C_i) \geq 2n/3$  but  $\text{cri}(C_{S(i)}) < 2n/3$ . Thus we have  $\text{cri}(C_{S(i)}) \in [n/3, 2n/3]$ .

**Correctness of the Reduction:** Fix  $C$  such that  $n/3 \leq \text{cri}(C) \leq 2n/3$ . Let

$$E' = \{(u, v) \in E \mid u \in \text{cri}(C), v \in V \setminus \text{cri}(C)\}.$$

We show that  $C$  contains every variable that appears in  $E'$ . If not, let  $e = (u, v) \in E'$  be a missing variable and suppose without loss of generality that  $u \in \text{cri}(C)$  and  $v \notin \text{cri}(C)$ . Since  $u \in \text{cri}(C)$ , by definition we know there exists a  $u$ -critical assignment  $\alpha_u$  such that  $C(\alpha_u) = 0$ . Let  $\alpha'_u$  be the same assignment but flipping  $x_{(u,v)}$ . Then by definition, we obtain a new assignment  $\alpha'_u$  that is  $v$ -critical. However, recall that  $v \notin \text{cri}(C)$ , which leads to a contradiction.

Thus, suppose that  $C$  is not obtained by an invalid derivation, then since  $\text{width}(C) < e$ , we know that  $|E'| < e$ , which means that  $\text{cri}(C)$  is a witness that the expansion of  $G$  is in fact less than  $e$ .

This finishes the proof.  $\square$

**Size Refuter.** After the PLS-membership of width refuter, we are ready to study the size refuter.

We consider Tseitin formulas where the underlying graph  $G = (V, E)$  is an expander. Recall from [Definition 6.6](#) that the expansion of  $G$ , denoted as  $e(G)$ , is the minimum number of edges between  $S$  and  $V \setminus S$  over every subset  $S \subseteq V$  such that  $|V|/3 \leq |S| \leq 2|V|/3$ . It is proved in [\[Sch97, BW01\]](#) that for every constant-degree expander  $G$  with  $e(G) \geq n$  and every odd-weighted function  $\tau : V \rightarrow \{0, 1\}$ , the tautology  $\text{Tseitin}(G, \tau)$  requires size- $2^{\Omega(n)}$  resolution proof.

Now, analogous to [Definition 6.8](#), we define the refuter problem for the size lower bounds, where the graph  $G$  is also given as an input, and a certificate for  $G$  not being an expander is also a valid output:

**Definition 6.11.** Let  $\text{REFUTER}(s(\text{Tseitin} \vdash_{\text{Res}} \perp) < 1.01^{n/d})$  denote the following problem. The input consists of an undirected connected  $d$ -regular graph  $G = (V, E)$  on  $n$  vertices, an odd-weighted assignment  $\tau : V \rightarrow \{0, 1\}$ , and a purported resolution refutation  $\Pi$  for  $\text{Tseitin}(G, \tau)$  that contains at most  $1.01^{n/d}$  clauses. A valid solution is either of the following:

- an index  $i$  such that the  $i$ -th node in  $\Pi$  is an invalid derivation, or
- a vertex set  $S \subseteq V$  such that  $|V|/3 \leq |S| \leq 2|V|/3$  and  $|E(S, V \setminus S)| < n$ .

Again, when we calculate the block-depth of reductions, we treat  $(G, \tau)$  as one input block.

**Theorem 6.12.** *Let  $G = (V, E)$  be a  $d$ -regular undirected connected graph and  $\tau : V \rightarrow \{0, 1\}$  be an odd-weighted function. Then, if  $e(G) \geq n$ , then  $\text{Tseitin}(G, \tau)$  requires resolution size  $\geq 1.01^{n/d}$ .*

*Moreover, there is a uniform decision tree reduction from  $\text{REFUTER}(s(\text{Tseitin} \vdash_{\text{Res}} \perp) < 1.01^{n/d})$  to  $\text{rwPHP}(\text{PLS})$  with block-depth 3.*

*Proof.* We follow the proof in [\[Sch97\]](#), which (also) uses a random restriction argument and a width lower bound. Our exposition about the random restrictions will be careful and slow (since this is relevant to our reduction to  $\text{rwPHP}(\text{PLS})$ ), but we will be sketchy about other parts.

Consider a random restriction as follows. Let  $t := n/10$ , pick  $t$  edges  $E' = \{e_1, e_2, \dots, e_t\}$  uniformly at random, and for each edge  $e_i$  assign a uniformly random bit to  $x_{e_i}$ . For an edge  $e = (x, y)$ , each time we assign  $x_e \leftarrow 0$ , we do nothing with the function  $\tau$ ; each time we assign  $x_e \leftarrow 1$ , we flip both  $\tau(x)$  and  $\tau(y)$ . After picking these  $t$  edges, we reduced the formula  $\text{Tseitin}(G, \tau)$  to the formula  $\text{Tseitin}(G', \tau')$ , where  $G'$  is the graph  $G$  with edges in  $E'$  removed, and  $\tau'$  is the assignment on vertices we obtained at the end. It is easy to see that  $e(G') \geq e(G) - t$ , hence by [Theorem 6.7](#), any resolution refutation for  $\text{Tseitin}(G', \tau')$  requires width  $\geq e(G) - t$ .

It would be helpful to rigorously define the space of random restrictions. Fix an ordering  $\prec$  (e.g., the lexicographic one) over the  $nd = 2|E|$  literals. A restriction is described by a sequence  $(i_0, i_1, \dots, i_{t-1})$  as follows. We first pick the  $i_0$ -th literal  $\ell_0$  according to  $\prec$  and set  $\ell_0 := 1$ . Now we are left with  $nd - 2$  literals (as both  $\ell_0$  and  $\bar{\ell}_0$  are set) and we pick the  $i_1$ -th literal  $\ell_1$  among them, according to  $\prec$ . After setting  $\ell_1 := 1$ , we are left with  $nd - 4$  literals and we pick the  $i_2$ -th one, and so on. Each sequence corresponds to a restriction that sets the values of  $t$  edges (but note that each restriction corresponds to  $t!$  such sequences). The space of random restrictions is denoted as

$$\mathcal{R} := [nd] \times [nd - 2] \times [nd - 4] \times \dots \times [nd - 2t + 2].$$

Let  $w := e(G) - t$  and fix a clause  $C$  of width  $\geq w$ . If we know that a restriction  $\rho$  does not kill  $C$ , then there is a more efficient way to describe  $\rho$  by a sequence  $(j_0, j_1, \dots, j_{t-1})$ , as follows. We first pick the  $j_0$ -th literal  $\ell_0$  among those  $nd - w$  ones not in  $C$ , according to  $\prec$ , and set  $\ell_0 := 1$ . After this round, there are at most  $nd - w - 1$  remaining literals not in  $C$ : If  $\bar{\ell}_0 \in C$  then there are exactly  $nd - w - 1$  such literals (i.e., excluding  $\ell_0$ ), otherwise there are  $nd - w - 2$  remaining literals (i.e., excluding  $\ell_0$  and  $\bar{\ell}_0$ ). Anyway, we use  $nd - w - 1$  as an upper bound on the number of literals that we can choose after the first round. In the next round, we choose the  $j_1$ -th literal  $\ell_1$  not in  $C$  according to  $\prec$ , set  $\ell_1 := 1$ , and now

there remains at most  $nd - w - 2$  literals. In the next round, we choose the  $j_2$ -th such literal, and so on. The space of “bad” restrictions that do not kill  $C$  is

$$\mathcal{BAD} := [nd - w] \times [nd - w - 1] \times \cdots \times [nd - w - t + 1].$$

Given any  $C$  and  $b \in \mathcal{BAD}$ , we can compute  $\text{seq}(C, b) \in \mathcal{R}$  as the “ $b$ -th bad restriction corresponding to  $C$ ”.<sup>29</sup> Given any clause  $C$  of width  $\geq w$  and any restriction  $\rho \in \mathcal{R}$  that does not kill  $C$ , we can compute an encoding  $\text{bad}(C, \rho) \in \mathcal{BAD}$  such that  $\text{seq}(C, \text{bad}(C, \rho)) = \rho$ . The following calculation corresponds to a “union bound” over  $L := 1.01^{n/d}$  clauses in the purported resolution proof:

$$\begin{aligned} L \cdot \frac{|\mathcal{BAD}|}{|\mathcal{R}|} &= L \cdot \prod_{i \in [t]} \frac{nd - e(G) + t - i}{nd - 2i} \\ &\leq L \cdot \left(1 - \frac{e(G) - 3t}{nd - 2t}\right)^t \\ &\leq 1.01^{n/d} \cdot \left(1 - \frac{7}{10d - 2}\right)^{n/10} \leq 1/2. \end{aligned}$$

The lower bound argument proceeds as follows. Let  $\Pi = (C_0, C_1, \dots, C_{L-1})$  be a purported size- $L$  resolution proof for  $\text{Tseitin}(G, \tau)$ . By the above union bound, there is a restriction  $\rho \in \mathcal{R}$  that kills every clause in  $\Pi$  with width  $\geq w$ . This restriction shrinks  $\Pi$  into  $\Pi|_\rho$  which is a width- $w$  resolution proof for  $\text{Tseitin}(G', \tau')$ , contradicting the width lower bound. Therefore, every resolution proof for  $\text{Tseitin}(G, \tau)$  requires more than  $1.01^{n/d}$  many clauses.

Finally, we describe the reduction from  $\text{REFUTER}(s(\text{Tseitin} \vdash_{\text{Res}} \perp) < 1.01^{n/d})$  to  $\text{rwPHP}(\text{PLS})$ :

( $f$ ) The function  $f : [L] \times \mathcal{BAD} \rightarrow \mathcal{R}$  is defined as  $f(i, b) := \text{seq}(C_i, b)$ .

( $I_\rho$ ) For every  $\rho \in \mathcal{SEQ}$ , we obtain a purported width- $w$  resolution proof  $\Pi|_\rho$  for  $\text{Tseitin}(G', \tau')$ . Every node in  $\Pi|_\rho$  can be computed in block-depth 1 from  $\Pi$ . Using [Theorem 6.9](#), we reduce the problem of finding an invalid derivation in  $\Pi|_\rho$  to an  $\text{ITER}$  instance  $I_\rho$ , where each node in  $I_\rho$  is computed in block-depth 3 from  $\Pi|_\rho$ .

( $g$ ) For every  $\rho \in \mathcal{SEQ}$  and every valid solution  $o$  of  $I_\rho$ , we can compute an index  $i \in [L]$  from  $o$  such that the  $i$ -th step in  $\Pi|_\rho$  is an illegal derivation. We let  $g_{\rho, o} := (i, \mathcal{BAD}(C_i, \rho))$ .

Suppose that  $(\rho, o)$  is any solution to the  $\text{rwPHP}(\text{PLS})$  instance defined above. Let  $i \in [L]$  be computed from  $o$  as above, then we claim that the  $i$ -th step of  $\Pi$  must be an illegal derivation. Indeed, since  $o$  is a solution of  $I_\rho$ , the  $i$ -th step of  $\Pi|_\rho$  must be illegal. On the other hand, if the  $i$ -th step of  $\Pi$  is not illegal, then  $C_i|_\rho$  is a clause of width  $\geq w$ , and thus

$$f(g_{\rho, o}) = \text{seq}(C_i, \mathcal{BAD}(C_i, \rho)) = \rho,$$

contradicting that  $(\rho, o)$  is a valid solution to the reduced  $\text{rwPHP}(\text{PLS})$  instance.  $\square$

## 6.4 Refuters for Random $k$ -CNFs

Finally, we show that resolution lower bounds for random  $k$ -CNFs can be refuted in  $\text{rwPHP}(\text{PLS})$ . More precisely, as in [\[CS88\]](#), we consider the distribution  $\mathcal{F}(k, n, m)$  over  $k$ -CNFs with  $n$  variables and  $m$  clauses where each clause is i.i.d. chosen from all  $\binom{n}{k} 2^k$  ordinary clauses of size  $k$  over the  $n$  variables. (A clause is *ordinary* if there is no variable  $x_i$  such that both  $x_i$  and  $\bar{x}_i$  occur in this clause.) Let  $c \geq 1, \varepsilon > 0$

<sup>29</sup>Note that some  $b \in \mathcal{BAD}$  might not correspond to a valid restriction. We can set  $\text{seq}(C, b)$  to be an arbitrary value.

be constants, and  $\{F_n\}_{n \in \mathbb{N}}$  be a family of  $k$ -CNFs, where each  $F_n$  is a  $k$ -CNF over  $n$  variables and  $cn$  clauses. In the search problem

$$\text{REFUTER}(s(F_n) < (1 + \varepsilon)^n),$$

we are given query access to a purported resolution refutation  $\Pi$  for  $F_n$  that contains at most  $(1 + \varepsilon)^n$  clauses, and our goal is to locate an invalid derivation in  $\Pi$ .

**Theorem 6.13.** *For every large enough positive integer  $k$  and  $c \geq 0.7 \cdot 2^k$ , there is a constant  $\varepsilon > 0$  such that the following holds. Let  $\{F_n\}_{n \in \mathbb{N}}$  be a sequence of random  $k$ -CNFs where each  $F_n$  is independently chosen according to the distribution  $\mathcal{F}(k, n, cn)$ . With probability 1, there is a non-uniform decision tree reduction of block-depth 2 from the problem  $\text{REFUTER}(s(F_n) < (1 + \varepsilon)^n)$  to  $\text{rwPHP(PLS)}$  that works for all large enough  $n$ .*

The unsatisfiability of  $\{F_n\}_{n \in \mathbb{N}}$  and the resolution lower bounds for  $\{F_n\}_{n \in \mathbb{N}}$  are already shown in the seminal work of Chvátal and Szemerédi [CS88]. We prove [Theorem 6.13](#) by formalizing their resolution lower bound proofs as decision tree reductions to  $\text{rwPHP(PLS)}$ .

The reason that our reduction in [Theorem 6.13](#) is non-uniform is very similar to that in [Section 6.1](#). First, it appears infeasible to decide if  $(1 + \varepsilon)^n$  is indeed a valid resolution size lower bound for the input formula  $F_n$ . Second, the proofs in [CS88] involve some objects that appear to be infeasible to compute given  $F_n$ ; however, these objects do not depend on the purported size- $(1 + \varepsilon)^n$  resolution refutation, thus can be hardwired in a non-uniform decision tree. It might be possible to obtain a “uniform version” of [Theorem 6.13](#) like what we did for Tseitin formulas ([Definition 6.8](#), [Remark 5](#), [Definition 6.11](#)), by completely formalizing [CS88] in bounded arithmetic. We choose not to do so because we believe that a non-uniform upper bound of  $\text{rwPHP(PLS)}$  already supports our claim that  $\text{rwPHP(PLS)}$  captures the complexity of proving *most* resolution lower bounds; dealing with extra details in [CS88] would only be distracting.

We assume familiarity with the (quite involved) proof in [CS88]. In particular, we need the following definitions and theorems:

- Fix a  $k$ -CNF  $F$  over  $n$  variables and  $cn$  clauses. Let  $X = \{x_1, x_2, \dots, x_n\}$  denote the set of variables of  $F$ . The “structure” of  $F$  can be described by a  $k$ -uniform (multi-)hypergraph  $H$  over the vertex set  $X$ , where each clause  $F_i$  of  $F$  corresponds to the hyperedge

$$E_i := \{x_j \in X : F_i \text{ contains } x_j \text{ or } \bar{x}_j\}.$$

- Let  $E'$  be a subset of hyperedges in  $H$ , the *boundary* of  $E'$  is the set of all vertices that belong to exactly one hyperedge in  $E'$ . We say that  $H$  has property  $P(a)$  if, for every  $m \leq an$ , every family of  $m$  edges has boundary size at least  $m/2$ .
- Let  $\mathcal{S}$  be a family of subsets of  $X$  (note that  $\mathcal{S}$  might be a multiset). A *system of distinct representatives* (SDR) of  $\mathcal{S}$  is a mapping from each  $S \in \mathcal{S}$  to an element in  $S$  such that different subsets in  $\mathcal{S}$  are mapped to different elements. Alternatively, consider the bipartite graph  $(\mathcal{S}, X)$  such that an edge between  $S \in \mathcal{S}$  and  $x_i \in X$  is drawn if and only if  $x_i \in S$ , then an SDR of  $\mathcal{S}$  is an  $\mathcal{S}$ -perfect matching of this bipartite graph (that is, every vertex in  $\mathcal{S}$  is matched).
- Let  $S$  be a subset of vertices in  $H$  of size  $s := \lfloor bn \rfloor$ . We say that  $S$  is *good* if there is a subset  $D$  of  $S$  with  $|S \setminus D| \leq (a/32)|S|$  such that every family of at most  $an$  edges has an SDR that is disjoint from  $D$ . We denote this subset as  $D(S)$ . We say that  $H$  has property  $Q(a, b)$  if a random size- $s$  subset  $S \subseteq X$  is good with probability at least  $1/2$ .
- [CS88, Lemma 3] showed that any hypergraph satisfying certain “sparsity” conditions will have properties  $P(a)$  and  $Q(a, b)$ . As a corollary ([CS88, Lemma 4]), for every large enough integers  $k$  and  $c \geq 0.7 \cdot 2^k$ , there are  $a, b > 0$  with  $b \leq a/8$  such that a random  $k$ -uniform hypergraph with

$n$  vertices and  $cn$  hyperedges has properties  $P(a)$  and  $Q(a, b)$  with probability  $\geq 1 - n^{-2}$  for large enough  $n$ .<sup>30</sup>

Now we outline the strategy of [CS88]. Let  $F = F_n$  be a  $k$ -CNF over  $n$  variables and  $cn$  clauses whose underlying hypergraph  $H$  satisfies  $P(a)$  and  $Q(a, b)$ . We first choose a “special pair”  $(S, \rho)$  where  $S$  is a random subset of  $s := \lfloor bn \rfloor$  vertices and  $\rho \in \{0, 1\}^{D(S)}$  is a uniformly random restriction on variables in  $D(S)$ . Then we use a *random restriction* argument to reduce the size lower bound to a *width lower bound*:

**Random restriction:** Let  $C$  be any clause in the purported resolution refutation for  $F$  such that the width of  $C$  is at least  $an/8$ . With probability  $1 - 2^{-\Omega(n)}$  over the choice of  $S$ , we have  $|\text{Vars}(C) \cap S| \geq as/16$ , where  $\text{Vars}(C)$  denotes the set of variables contained in  $C$ . Since  $|S \setminus D(S)| \leq as/32$ , it follows that  $|\text{Vars}(C) \cap D(S)| \geq as/32$ , hence the probability over  $\rho$  that  $C$  is not killed by  $\rho$  is at most  $2^{-as/32}$ . A union bound over all  $C \in \Pi$  implies that with high probability over  $(S, \rho)$ , every clause of width  $\geq an/8$  in  $\Pi$  is killed by  $\rho$ .

**Width lower bound:** Now we are left with a purported resolution refutation  $\Pi|_\rho$  of width less than  $an/8$  for the statement  $F|_\rho$ . For a clause  $C \in \Pi$ , let  $\mu(C)$  denote the minimum number of clauses from  $F$  that logically implies  $C$  under  $\rho$ . (That is, for any assignment extending  $\rho$ , if all these clauses are satisfied, then  $C$  is also satisfied.) Every subset of  $\leq an/2$  clauses  $F' \subseteq F$  can be satisfied by some assignment on  $X \setminus D(S)$  (indeed, we can simply choose an SDR for  $F'$  that is disjoint from  $D(S)$ , and fix this SDR), hence  $\mu(\perp) > an/2$ . On the other hand, every clause in  $F$  has  $\mu$  value at most 1. Let  $C' \in \Pi|_\rho$  be the first clause in  $\Pi|_\rho$  such that  $\mu(C') > an/2$  (recall that  $\perp$  is the last clause in  $\Pi|_\rho$ ). One can use a classical argument to show that  $an/2 < \mu(C') \leq an$ , i.e., the smallest subset of clauses  $F' \subseteq F$  that logically implies  $C'$  has size between  $an/2$  and  $an$ . Since  $H$  satisfies  $P(a)$  and  $|F'| \leq an$ , the boundary of  $F'$  contains at least  $|F'|/2 \geq an/4$  variables. It can be shown that  $C'$  contains every variable in the boundary of  $F'$  but not in  $S$ , hence  $w(C') \geq an/8$ , a contradiction.

Now we are ready to prove [Theorem 6.13](#).

*Proof of Theorem 6.13.* Let  $F = F_n$  be the random  $k$ -CNF and  $H = H_n$  be the underlying hypergraph for  $F$ . Let  $a, b > 0$  be constants that arise from [CS88, Lemma 4], we assume that  $H$  has properties  $P(a)$  and  $Q(a, b)$  (this assumption will be justified at the end of the proof). Our reduction needs the following non-uniform advice  $\{S_i\}, \{D_i\}, \{\mathcal{R}_{i,\rho}\}$  (of course, they only depend on  $F$  and is independent of the purported resolution refutation):

- A list of subsets  $S \subseteq X$  with size  $s := \lfloor bn \rfloor$  that are good. Since  $H$  has property  $Q(a, b)$ , there are at least  $N_{\text{good}} := \binom{n}{s}/2$  such subsets and we only need to encode the first  $N_{\text{good}}$  ones. For each  $i \in [N_{\text{good}}]$ , denote the  $i$ -th good subset as  $S_i$ , we also need the subset  $D_i \subseteq S_i$  of size  $\geq (1 - a/32)s$  such that every family of at most  $an$  edges has an SDR disjoint from  $D_i$ .
- For each index  $i$  and each restriction  $\rho \in \{0, 1\}^{D_i}$ , we compute the subformula  $F|_\rho$ . The above width lower bound argument (along with properties  $P(a)$  and  $Q(a, b)$ ) implies that  $F|_\rho$  requires resolution width  $> an/8$ . Invoking [Theorem 6.1](#), we obtain a non-uniform decision tree reduction from  $\text{REFUTER}(w(F|_\rho) \leq an/8)$  to PLS with block-depth 2, which we denote as  $\mathcal{R}_{i,\rho}$ .

Let  $\Pi$  be a purported resolution refutation for  $F$  consisting of at most  $L := (1 + \varepsilon)^n$  clauses. Now we describe our reduction from  $\text{REFUTER}(s(F_n) \leq (1 + \varepsilon)^n)$  to  $\text{rwPHP}(\text{PLS})$ :

<sup>30</sup>If this probability is at least  $1 - n^{-2}$ , then we can argue that with probability 1 over an infinite family of random  $k$ -CNFs, our reduction to  $\text{rwPHP}(\text{PLS})$  is correct on all but finitely many input lengths; see the end of the proof of [Theorem 6.13](#). Although [CS88] only claimed a probability of  $1 - o(1)$ , their proof actually shows a probability of  $1 - n^{-\Omega(k)}$  where the big  $\Omega$  hides some absolute constant. This is at least  $1 - n^{-2}$  when  $k$  is large enough; we suspect that our results can be extended to all  $k \geq 3$  via a more careful argument.

(*f*) The function  $f$  takes as inputs  $i \in [L]$ ,  $type \in \{0, 1\}$ , and  $w \in [\binom{n}{s} \cdot 2^{(1-a/32)s} \cdot 2^{-c'n}]$ , where  $c' > 0$  is a small enough constant depending on  $a$  and  $b$ . Essentially, it treats  $(i, type, w)$  as the compression of a bad “special pair”  $(S, \rho)$  (where  $S$  is a good size- $s$  subset and  $\rho$  is an assignment over  $D(S)$ ) and decompresses it. We start by checking that  $w(C_i) \geq an/8$ ; if this is not the case then  $f$  outputs  $\perp$ . Next:

- If  $type = 0$ , then this means  $|\text{Vars}(C_i) \cap S| < as/16$ . Recall that if  $|S| = s$  is chosen uniformly at random, then the probability that  $|\text{Vars}(C_i) \cap S| < as/16 \leq 0.5 \cdot |C_i|s/n$  should be at most  $2^{-c'n}$  for some small enough constant  $c' > 0$ . Hence,  $(S, \rho)$  can be compressed into  $(\log \binom{n}{s} - c'n) + |D(S)|$  bits. We treat  $w$  as this compression and recover  $(S, \rho)$  from  $w$ .
- If  $type = 1$ , then  $|\text{Vars}(C_i) \cap S| \geq as/16$  but  $C_i$  is not killed under  $\rho$ . In this case, the values of  $\rho$  over  $\text{Vars}(C_i) \cap D(S)$  can be inferred from  $C_i$ . Since  $|\text{Vars}(C_i) \cap D(S)| \geq as/16 - as/32 = as/32$ , this provides us a way to compress  $(S, \rho)$  into  $\log \binom{n}{s} + (|D(S)| - as/32) \leq \log \binom{n}{s} + |D(S)| - c'n$  bits. Again, we treat  $w$  as this compression and recover  $(S, \rho)$  from  $w$ .

Now that we obtained  $(S, \rho)$ , we can find an index  $j \in [N_{\text{good}}]$  such that  $S = S_j$  (using non-uniformity). If such  $j$  does not exist, then  $f$  outputs  $\perp$ ; otherwise  $f$  outputs  $(j, \rho)$ .

Hence we have  $f : [L] \times \{0, 1\} \times [\binom{n}{s} \cdot 2^{(1-a/32)s} \cdot 2^{-c'n}] \rightarrow [N_{\text{good}}] \times \{0, 1\}^{(1-a/32)s}$ . (If  $f$  outputs  $\perp$  then we can assume that it outputs a default value, say  $(0, 0^{(1-a/32)s})$ , instead.) Recall that  $L = (1 + \varepsilon)^n$  and  $N_{\text{good}} = \binom{n}{s}/2$ , which means if  $\varepsilon > 0$  is small enough then

$$\frac{2L \cdot \binom{n}{s} \cdot 2^{(1-a/32)s} \cdot 2^{-c'n}}{N_{\text{good}} \cdot 2^{(1-a/32)s}} \leq 2^{-\Omega(n)} \ll 1, \quad (8)$$

hence  $f$  is indeed shrinking. Given an input  $(i, type, w)$ , its  $f$  value can be computed by a non-uniform decision tree of block-depth 1.

( $I_{j,\rho}$ ) Given  $j \in [N_{\text{good}}]$  and  $\rho \in \{0, 1\}^{(1-a/32)s}$ , we compute a PLS instance  $I_{j,\rho}$  as follows. Abusing notation, we also use  $\rho$  to denote the restriction that equals to  $\rho$  on  $D_j$  and does not restrict any variable outside  $D_j$ . Let  $\Pi|_\rho$  denote the restriction of  $\Pi$  over  $\rho$ , then each clause of  $\Pi|_\rho$  can be computed in block-depth 1 from  $\Pi$ . Then we apply the reduction  $\mathcal{R}_{i,\rho}$  on  $\Pi|_\rho$  to obtain the PLS instance  $I_{j,\rho}$ .

(*g*) Let  $j \in [N_{\text{good}}]$  and  $\rho \in \{0, 1\}^{(1-a/32)s}$ . Given a valid solution  $o$  of  $I_{j,\rho}$ , we can compute an index  $i \in [L]$  from  $o$  such that the  $i$ -th step in  $\Pi|_\rho$  is an illegal derivation. As in the definition of  $f$ , we can (assume  $w(C_i) \geq an/8$  and) compress  $(j, \rho)$  as  $(i, type, w)$ ; then we set  $g_{(j,\rho),o} = (i, type, w)$ . If  $f(i, type, w) \neq (j, \rho)$ , then it must be the case that the  $i$ -th step in  $\Pi$  is already incorrect (instead of the case that  $w(C_i)$  is too large).

The above reduction is correct as long as  $H$  has properties  $P(a)$  and  $Q(a, b)$ , and its block-depth is 2.

It remains to show that our reduction is correct with probability 1. In fact, for each  $N \geq 1$ , the probability that for every  $n \geq N$ ,  $H_n$  has properties  $P(a)$  and  $Q(a, b)$  is at least

$$\prod_{n \geq N} (1 - n^{-2}) = \frac{N-1}{N}.$$

It follows that with probability 1 over the family  $\{F_n\}_{n \in \mathbb{N}}$ , all but finitely many  $H_n$  has properties  $P(a)$  and  $Q(a, b)$ . In this case, our reduction will be correct on all but finitely many input lengths.  $\square$

## 6.5 Open Problems: What We *Failed* to Formalize

One interesting problem left open by this work is whether the general size-width trade-offs in [BW01] can be proved in  $\text{rwPHP}(\text{PLS})$ . Ben-Sasson and Wigderson showed that for any unsatisfiable  $k$ -CNF  $F$ , if  $F$  requires resolution width  $w$  to refute, then  $F$  also requires resolution size  $2^{\Omega(w-k)^2/n}$  to refute. This naturally leads to the following conjecture:

**Conjecture 6.14** (Informal). *Let  $F$  be an unsatisfiable  $k$ -CNF with resolution width  $> w_F$  and let  $s_F := 2^{\Omega(w_F-k)^2/n}$ . Let  $\mathcal{P}$  denote the problem  $\text{REFUTER}(w(F \vdash_{\text{Res}} \perp) \leq w_F)$ , then there is an efficient decision-tree reduction from  $\text{REFUTER}(s(F \vdash_{\text{Res}} \perp) \leq s_F)$  to  $\text{rwPHP}(\mathcal{P})$ . In particular, there is always an efficient non-uniform decision tree reduction from  $\text{REFUTER}(s(F \vdash_{\text{Res}} \perp) \leq s_F)$  to  $\text{rwPHP}(\text{PLS})$ .*

Roughly speaking, one obstacle against proving **Conjecture 6.14** is that the averaging argument used in the proof of [BW01, Theorem 3.5] seems to rely on “APC<sub>2</sub>-style” [Jeř09] approximate counting: one needs to estimate the number of “fat” clauses up to an  $(1 + \varepsilon)$ -multiplicative factor. Therefore, we have been unable to formalize the proof of [BW01, Theorem 3.5] in  $\text{T}_2^1 + \text{dwPHP}(\text{PV})$  where only “APC<sub>1</sub>-style” [Jeř07a] approximate counting is available.

We also leave open the complexity of proving resolution lower bounds by combining monotone circuit lower bounds [Raz85, AB87, Hak95] with feasible interpolation [Raz95b, Kra97, Pud97]. To formalize Razborov’s approximation method [Raz85], it seems that we need to iteratively define exponentially many set families (one for each node in the resolution proof) and apply the sunflower lemma [ER60, ALWZ21] to each of them. It is unclear to us how to formalize such arguments in  $\text{T}_2^1 + \text{dwPHP}(\text{PV})$ . (See also [GGKS20] who used lifting techniques to prove monotone circuit lower bounds and resolution lower bounds.)

We showed in **Corollary 6.2** that the refuter problem for every true resolution width lower bound is  $\text{PLS}$ -complete under non-uniform reductions. It would be very interesting to see whether the size lower bound analog holds or not. We propose the following conjecture (which is stronger than the non-uniform version of **Conjecture 6.14**):

**Conjecture 6.15** (Informal). *Let  $F$  be an unsatisfiable CNF that requires resolution size  $\geq s_F$  to refute. Then the problem  $\text{REFUTER}(s(F \vdash_{\text{Res}} \perp) < s_F)$  is  $\text{rwPHP}(\text{PLS})$ -complete under non-uniform decision tree reductions.*

(Note that the *average-case* version of **Conjecture 6.15**, where  $F$  is a random  $k$ -CNF and  $s_F = 2^{\Omega(n)}$ , is already proved in **Section 6.3**, by formalizing the resolution size lower bounds of [CS88].)

We end this subsection by mentioning a subtle technical issue in our proofs. There are two natural properties in the completeness of resolution (i.e., resolution can prove every true statement within size  $2^n$ ): the proof does not require weakening, and it avoids producing duplicate clauses. However, in our current  $\text{PLS}$ -hardness of refuting resolution width lower bounds and  $\text{rwPHP}(\text{PLS})$ -hardness of refuting resolution size lower bounds, the resolution proofs produced in our reduction rely on both weakening rules and duplicated clauses. This raises an open question: What is the complexity of the corresponding refuter problems if the proofs are restricted from using either weakening rules or duplicate clauses?

## 7 Applications

### 7.1 Proof Complexity of Proof Complexity Lower Bounds

In this subsection, we translate our TFNP upper bounds for the refuter problems into *proof complexity upper bounds* for *proof complexity lower bounds*, showing that resolution lower bounds can actually be proved in weak proof systems! In particular, we show that low-width resolution (itself) can prove lower bounds on resolution width (**Theorem 7.1**), while low-width *random resolution* (as defined in [BKT14,

PT19]) can prove resolution size lower bounds (Theorem 7.3).<sup>31</sup> This stands in stark contrast to the results proven in [AM20, Gar19, dRGN<sup>+</sup>21] that resolution cannot prove size lower bounds against itself.

**Formalization of proof complexity lower bounds as CNFs.** Suppose a family of formulas  $\mathcal{F} = \{F_n\}$  does not have a width  $w_F$  resolution refutation (i.e.,  $w(\mathcal{F} \vdash_{\text{Res}} \perp) > w_F$ ). Then we can transform the refuter problem  $\text{REFUTER}(w(\mathcal{F} \vdash_{\text{Res}} \perp) \leq w_F)$  into a family of unsatisfiable CNFs  $\mathcal{F}_{\text{wLB}}^w$  using via false clause search problem (Equation 2). That is, an unsatisfiable CNF  $F_{\text{wLB}}^w$  in the family  $\mathcal{F}_{\text{wLB}}^w$  is defined as follows:

- The input of  $F_{\text{wLB}}^w$  is a purported length- $L$  resolution refutation for  $F_n$  represented as a list of nodes  $C_0, C_1, \dots, C_{L-1}$  and each node  $C_i$  can be encoded in  $O(w_F \log n)$  bits.
- Each potential solution  $sol$  of the refuter problem can be verified by a decision tree of block-depth 3, hence they can be turned into a CNF  $C_{sol}$  of width  $O(w_F \log n)$ .  $F_{\text{wLB}}^w$  is simply the conjunction of these CNFs.

We can similarly transform a resolution *size* lower bound  $s(\mathcal{F} \vdash_{\text{Res}} \perp) > L$  into a family of unsatisfiable CNFs  $\mathcal{F}_{\text{sLB}}^L$  via the refuter problem  $\text{REFUTER}(s(\mathcal{F} \vdash_{\text{Res}} \perp) \leq L)$ . The only difference is that each node consists of an (unbounded-width) clause and thus is encoded in  $O(n + \log L)$  bits.

It is easily seen that  $\mathcal{F}_{\text{wLB}}^w$  are CNFs of width  $O(w_F \log n)$  and  $\mathcal{F}_{\text{sLB}}^L$  are CNFs of width  $O(n + \log L)$ . (When  $L = 2^{n^{\Omega(1)}}$ , these width parameters are  $\text{polylog}(L)$  and can be thought of as “efficient”.)

*Remark 6* (Comparison with previous formalizations). Similar formalizations of resolution lower bound statements have also appeared in [AM20, Gar19, dRGN<sup>+</sup>21]. The biggest difference between these formalizations is that in [dRGN<sup>+</sup>21], the predecessors of each node are represented in binary and as  $O(\log N)$  bits; while in [AM20, Gar19], the predecessors are represented in unary and we have tables  $L[i, j]$  and  $R[i, j]$  denoting whether node  $j$  is a predecessor of node  $i$ . Note that in the unary representation, it requires an axiom of width  $L$  to express that every node  $u$  has at least one predecessor  $L[u]$  and at least one predecessor  $R[u]$ . Thus it is impossible to prove resolution width lower bounds in resolution width  $O(w \log N) \ll L$ . Therefore, we choose to use the binary formalization as in [dRGN<sup>+</sup>21].

The formalization in [dRGN<sup>+</sup>21] allows *disabled* nodes in the resolution proof. Our proof complexity upper bounds hold regardless of whether such nodes are allowed in the formalization.

**Low-width resolution can prove resolution width lower bounds.** First, we show that:

**Theorem 7.1.** *For every family of unsatisfiable CNFs  $\mathcal{F}$ , if  $w(\mathcal{F} \vdash_{\text{Res}} \perp) > w_F$ , then  $w(\mathcal{F}_{\text{wLB}}^w \vdash_{\text{Res}} \perp) \leq O(w_F \log N)$ .*

Theorem 7.1 follows from the proof of Theorem 6.1 and Theorem 3.4: since the refuter problem corresponding to resolution width lower bounds can be solved in PLS and the totality of PLS can be proved in low resolution width, it follows that resolution width lower bounds themselves can be proved in low resolution width. For the sake of intuition, we also present an equivalent but more direct proof using *Prover-Delayer games* [Pud00]. The necessary backgrounds on Prover-Delayer games are presented in Section C.1.

*Proof.* It suffices to construct a Prover strategy with memory size  $O(w_F \log N)$  in the Prover-Delayer game for  $\mathcal{F}_{\text{wLB}}^w$ . The Prover starts by querying the last node in the purported resolution proof, which should contain the empty clause  $\perp$ . The Prover maintains the invariant that she is always at some (not disabled) clause  $C_i$  such that  $w(\mathcal{F} \vdash C_i) > w_F$ , i.e., it requires resolution width  $> w_F$  to derive  $C_i$  from the axioms. Each time the Prover is at some clause  $C_i$ :

<sup>31</sup>More precisely, we use Theorem 6.1 to show that low-width resolution can prove *every* resolution width lower bound that is true, and use Theorem 6.13 to show that low-width random resolution can prove *most* resolution size lower bounds.

- Suppose  $C_i$  is *resolved* from the clauses  $C_j, C_k$ . Then the Prover queries  $C_j$  and  $C_k$ ; if  $j \geq i, k \geq i$ , or the derivation from  $(C_j, C_k)$  to  $C_i$  is invalid, then she wins the game. Otherwise, since the widths of  $C_j$  and  $C_k$  are at most  $w_F$  (recall that this is guaranteed syntactically by only allocating  $w_F$  variables to each clause), one of  $C_j, C_k$  must require  $> w_F$  width to derive. Suppose it is  $C_j$ ; that is,  $w(\mathcal{F} \vdash C_j) > w_F$ . Then the Prover forgets  $C_i$  and  $C_k$  and only remembers  $C_j$ .
- Suppose  $C_i$  is a *weakening* of a clause  $C_j$ . The Prover queries  $C_j$ ; if  $j \geq i$  or the weakening from  $C_j$  to  $C_i$  is invalid, then she wins the game. Otherwise, it must be the case that  $w(\mathcal{F} \vdash C_j) > w_F$ . Then the Prover forgets  $C_i$  and only remembers  $C_j$ .

Since the index  $i$  is always decreasing, the Prover is guaranteed to win the game. The Prover only needs to memorize  $O(1)$  resolution nodes, i.e.,  $O(w_F \log N)$  bits.  $\square$

**Low-width random resolution can prove resolution size lower bounds.** We first define the random resolution system (denoted as rRes):

**Definition 7.2** ([BKT14, PT19]). An  $\varepsilon$ -random resolution refutation of an unsatisfiable formula  $F$  is a distribution  $\mathcal{D}$  supported on pairs  $(\Pi, B)$ , such that

1. each  $B$  is a CNF formula over the variables of  $F$ ,
2.  $\Pi$  is a resolution refutation of  $F \wedge B$ , and
3. for any assignment  $x \in \{0, 1\}^n$ ,  $\Pr_{(\Pi, B) \sim \mathcal{D}}[B(x) = 1] \geq 1 - \varepsilon$ .

The *size*  $s(F \vdash_{\text{rRes}} \perp)$ , and *width*  $w(F \vdash_{\text{rRes}} \perp)$  of a random resolution refutation  $\mathcal{D}$  for  $F$  are the maximum size and width of a proof  $\Pi$  in the support of  $\mathcal{D}$ , respectively.

We remark that random resolution is not a standard (i.e., Cook–Reckhow) proof system since the distribution  $\mathcal{D}$  might potentially require exponentially many bits to describe and it is also unclear how to verify [Item 3](#) above. (In fact, random resolution cannot be simulated by a Cook–Reckhow proof system unless  $\text{P} = \text{NP}$  [PT19, Proposition 3.3].) On the other hand, strong lower bounds on both width and size are known for random resolution [PT19], suggesting that it may be classified as a “weak” proof system.

**Theorem 7.3.** *For every  $k \geq 3$  and  $c \geq 0.7 \cdot 2^k$ , there exists some  $\varepsilon > 0$  such that the following holds. Let  $F$  be a random  $k$ -CNF formula chosen from the distribution  $\mathcal{F}(k, n, cn)$ ,  $L := (1 + \varepsilon)^n$ , and  $\mathcal{F}_{\text{sLB}}^L(F)$  be the CNF formula encoding the lower bound that  $F$  requires size- $L$  resolution refutation. With probability tending to 1 (when  $n \rightarrow \infty$ ) over  $F$ ,  $\mathcal{F}_{\text{sLB}}^L(F)$  admits a  $\text{poly}(n)$ -width  $\gamma$ -random resolution refutation with  $\gamma := 2^{-\Omega(n)}$ .*

Similarly, [Theorem 7.3](#) is a corollary of [Theorem 6.13](#): if a search problem reduces to rwPHP(PLS), then it also *randomly* reduces to PLS, and such a random reduction can be translated into a random resolution refutation. Nevertheless, for the sake of intuition, we present an (equivalent) proof that directly constructs the random resolution refutation  $(\Pi, B)$ .

*Proof Sketch.* We assume familiarity with the proofs in [Section 6.4](#). We use the parameters  $a, b$  from [CS88, Lemma 4], and denote  $s := \lfloor bn \rfloor$ . We assume that the properties  $P(a)$  and  $Q(a, b)$  holds for  $F$ ; by [CS88, Lemma 4], this is true with high probability over  $F \leftarrow \mathcal{F}(k, n, cn)$ .

Recall that the variables in  $\mathcal{F}_{\text{sLB}}^L(F)$  encode a length- $L$  resolution refutation  $C_0, \dots, C_{L-1}$  of  $F$ , where  $L := (1 + \varepsilon)^n$ . Let  $S_0, S_1, \dots, S_{N_{\text{good}}-1}$  denote the first  $N_{\text{good}} := \binom{n}{s}/2$  good size- $s$  subsets. For each  $j \in [N_{\text{good}}]$ , also let  $D_j$  denote any subset of  $S_j$  of size  $\geq (1 - a/32)s$  such that every family of at most  $an$  edges has an SDR disjoint from  $D_j$ . To sample a pair  $(\Pi, B)$ :

1. We first pick a random  $j \in [N_{\text{good}}]$  and then pick a string  $\rho \leftarrow \{0, 1\}^{D_j}$ . We also treat  $\rho$  as a restriction that fixes every variable in  $D_j$  and leaves everything else unchanged.

2. For each  $i \in [L]$ , let  $B_i$  be the decision tree verifying that either  $w(C_i) < an/8$  or  $\rho$  kills  $C_i$ . Note that  $B_i$  only depends on the clause  $C_i$ , which can be encoded in  $\text{poly}(n)$  bits. Let  $B := \bigwedge_{i \in [L]} B_i$ , then  $B$  is a  $\text{poly}(n)$ -width CNF.

Moreover, following the same calculation as [Equation 8](#), we can show that for any assignment  $x$  to the variables in  $\mathcal{F}_{\text{sLB}}^L(F)$  (i.e.,  $x$  encodes a purported resolution refutation of  $F$ ), the probability over  $B$  (i.e., over  $j$  and  $\rho$ ) that  $B(x) = 1$  is at least  $1 - 2^{-\Omega(n)}$ .

3. It remains to argue that there always exists a  $\text{poly}(n)$ -width resolution refutation  $\Pi$  for  $\mathcal{F}_{\text{sLB}}^L(F) \wedge B$ . We can use a similar Prover's strategy as described in [Theorem 7.1](#). Recall that for any clause  $C$ ,  $\mu(C)$  denotes the minimum number of clauses from  $F$  that logically implies  $C$  under  $\rho$ , and that  $\mu(\perp) > an/2$ . The Prover starts from  $C_{L-1} = \perp$  and maintains the invariant that she is always at some clause  $C_i$  where  $\mu(C_i) > an/2$ . In addition, when the Prover is at some clause  $C_i$ , she also ensures that  $B_i$  is satisfied. At some stage, she will encounter some  $C_i$  that is resolved from  $C_j, C_k$ , such that  $\mu(C_i) > an/2$  and  $\mu(C_j), \mu(C_k) < an/2$ . But due to the width lower bound in [Section 6.4](#), this will imply that either the  $i$ -th derivation is invalid, or that  $B_i$  is violated.

It is easy to check that this Prover strategy only requires  $\text{poly}(n)$  memory.  $\square$

## 7.2 Complexity of Black-Box TFNP Separations

In this subsection, we introduce a new type of refuter problems —  $\text{TFNP}^{\text{dt}}$  *refuter* — which corresponds to the “complexity” of proving black-box  $\text{TFNP}^{\text{dt}}$  separations. We present the definition and several basic properties of them in [Section 7.2.1](#). In [Section 7.2.2](#), we relate the  $\text{TFNP}^{\text{dt}}$  refuter to the resolution width refuter ([Lemma 7.10](#)). Combining this with our results on resolution width refuter for EPHP and Tseitin, we characterize the “complexity” of separating PPA and PPP from PLS in the black-box setting by the class PLS itself.

**Notations.** For two  $\text{TFNP}^{\text{dt}}$  problems  $\mathcal{P}, \mathcal{Q}$ , we write  $\mathcal{P} \leq_m \mathcal{Q}$  if there is a *many-one* reduction from  $\mathcal{P}$  to  $\mathcal{Q}$ ; if the reduction is also uniform, we write  $\mathcal{P} \leq_m^U \mathcal{Q}$ .

### 7.2.1 Black-Box TFNP Refuters and its Properties

We start by providing a formal definition of the  $\text{TFNP}^{\text{dt}}$  refuter problems. Roughly speaking, in the problem  $\text{REFUTER}_{d,M}(\mathcal{P} \rightarrow \mathcal{Q})$ , we are given a shallow decision tree that claims to reduce  $\mathcal{P}$  to  $\mathcal{Q}$ , and our goal is to find a witness that this shallow decision tree is incorrect.

#### Problem $\text{REFUTER}_{d,M}(\mathcal{P} \rightarrow \mathcal{Q})$

Parameters: Two  $\text{TFNP}^{\text{dt}}$  problems  $\mathcal{P} = \{P_N\}, \mathcal{Q} = \{Q_N\}$  and two functions  $d := d(N), M := M(N)$  such that there is no depth- $d(N)$  decision tree reduction from  $P_N$  to  $Q_{M(N)}$  for any  $N$ .

Input: A purported depth- $d$  decision tree reduction  $(f_i, g_o)_{i \in M, o \in O_Q}$  from  $P_N = \{0, 1\}^N \times O_P$  to  $Q_M = \{0, 1\}^M \times O_Q$ .

Output: A pair  $(\rho, o^*)$ , where

- $\rho \in \{0, 1, *\}^N$  is a partial assignment encoded by specifying the locations and the values of all non- $*$  bits;
- $o^* \in O_Q$  is a solution of the problem  $Q_M$ .

The pair  $(\rho, o^*)$  satisfy that for any input  $x \in \{0, 1\}^N$  consistent with  $\rho$ ,

1.  $(f(x), o^*) \in \mathcal{Q}$  and  $(x, g_{o^*}(x)) \notin \mathcal{P}$ ;

2. only bits specified in  $\rho$  are ever queried when calculating  $g_{o^*}(x)$  and verifying  $(f(x), o^*) \in \mathcal{Q}$  and  $(x, g_{o^*}(x)) \notin \mathcal{P}$ .

In this section, we only consider refuting *low-depth* many-one decision tree reduction. Thus, we always assume the functions  $d(N)$  and  $\log M(N)$  are poly-logarithmic in  $N$  when we write “for any  $d, M$ ”. We also assume  $M(N) \geq N$ , so we will not consider reductions that are too weak. Note that it is necessary to have  $o^*$  as part of the solution for this problem to be in  $\text{TFNP}^{\text{dt}}$ ; otherwise, it might take too many queries to the input (reduction) to find  $o^*$ , which is used to refute the reduction later.

We call a  $\text{TFNP}^{\text{dt}}$  problem  $\mathcal{R}$  *syntactical* if all the decision trees ( $T_o$ ) for verifying the solution can be replaced by a single polynomial-time oracle Turing machine.<sup>32</sup> Since we mostly care about the black-box separations between *syntactical*  $\text{TFNP}$  subclasses, we assume all the  $\text{TFNP}^{\text{dt}}$  problems in this section are syntactical.

We now present several basic properties regarding the  $\text{TFNP}^{\text{dt}}$  separation refuter. First, a weaker reduction, which has a lower depth or smaller instance size, is easier to refute. The proof trivially follows from the definition (where item 2 needs the problem  $\mathcal{Q}$  to be paddable).

**Lemma 7.4.** 1. If  $d_1 \leq d_2$ , then  $\text{REFUTER}_{d_1, M}(\mathcal{P} \rightarrow \mathcal{Q}) \leq_m^U \text{REFUTER}_{d_2, M}(\mathcal{P} \rightarrow \mathcal{Q})$ .

2. If  $M_1 \leq M_2$ , then  $\text{REFUTER}_{d, M_1}(\mathcal{P} \rightarrow \mathcal{Q}) \leq_m^U \text{REFUTER}_{d, M_2}(\mathcal{P} \rightarrow \mathcal{Q})$ .

Even in the easiest parameter settings, i.e.,  $d = 0$ ,  $M(N) = N$ , the refuter problem  $\text{REFUTER}_{0, N}(\mathcal{P} \rightarrow \mathcal{Q})$  is least as hard as  $\mathcal{Q}$  itself, because a valid solution of  $\mathcal{Q}$  is always required to witness a mistake given by the input reduction.

**Lemma 7.5.**  $\mathcal{Q} \leq_m^U \text{REFUTER}_{0, N}(\mathcal{P} \rightarrow \mathcal{Q})$ .

*Proof.* Let  $y \in \{0, 1\}^N$  be an instance of problem  $Q_N \in \{0, 1\}^N \times O_Q$  and let  $o_P$  be an arbitrary fixed solution of problem  $P_N$ . We construct a trivial depth-0 reduction  $(f_i, g_o)_{i \in N, o \in O_Q}$ , where

$$f_i(x) = y_i, \forall i \in N; g_o(x) = o_P, \forall o \in O_Q.$$

Consider such reduction  $(f_i, g_o)$  as an instance of  $\text{REFUTER}_{0, N}(\mathcal{P} \rightarrow \mathcal{Q})$ , and let  $(\rho, o^*)$  be any solution of it. By definition,  $o^*$  is a valid solution of instance  $y$ . Moreover, our reduction is uniform, though  $(f_i, g_o)$  is not.  $\square$

Finally, we present two useful lemmas, which state that it is easier to refute a reduction when the *difficulty gap* between these two problems becomes larger.

**Lemma 7.6.** If  $\mathcal{P} \leq_m^U \mathcal{S}$ , and  $d_1(N), \log M_1(N) = \text{polylog}(N)$ , then

$$\text{REFUTER}_{d_1, M_1}(\mathcal{S} \rightarrow \mathcal{Q}) \leq_m^U \text{REFUTER}_{d_2, M_2}(\mathcal{P} \rightarrow \mathcal{Q}),$$

for some  $d_2(N), \log M_2(N) = \text{polylog}(N)$ .

*Proof.* Given a depth- $d_1$  reduction  $(f_i, g_o)_{i \in M_1, o \in O_Q}$  from  $S_N$  to  $Q_{M_1}$ , we compose it with any (uniform) low-depth reduction  $(h_i, l_o)_{i \in N, o \in O_S}$  from  $P_{N'}$  to  $S_N$ . Now we get a depth- $d'$  reduction  $(f'_i, g'_o)$  from  $P_{N'}$  to  $Q_{M_1}$  with

$$f'_i(x) = f_i(h(x)), \forall i \in [M_1]; \quad g'_o(x) = l_s(x), s := g_o(h(x)), \forall o \in O_Q.$$

Let  $d_2(N') := d'$ ,  $M_2(N') := M_1$ , and it is easy to verify that  $d_2(N'), \log M_2(N') = \text{polylog}(N')$ .

Consider a pair  $(\rho_P, o^*)$  that refutes  $(f'_i, g'_o)$ . Let  $x$  be any input of  $P_{N'}$  that is consistent with  $\rho_P$  and define  $y := h(x)$ . We show how to construct a partial assignment  $\rho_S$  consistent with  $y$  such that  $(\rho_S, o^*)$  refutes  $(f_i, g_o)$ . We start with setting  $\rho_S$  to all  $*$  strings, and then execute the process of

<sup>32</sup>A *syntactical*  $\text{TFNP}^{\text{dt}}$  problem is essentially a type-2  $\text{TFNP}$  ( $\text{TFNP}^2$ ) problem, see [BCE<sup>+</sup>98].

**S1:** calculating  $g_{o^*}(y)$ ;      **S2:** verifying  $(f(y), o^*) \in \mathcal{Q}$ ;      **S3:** verifying  $(y, g_{o^*}(y)) \notin \mathcal{S}$ .

During the above process, if  $y_i$  is queried and  $y_i$  has not been specified by  $\rho_S$ , we will execute  $h_i(x)$  to calculate  $y_i$  and then store its value in  $\rho_S$ .

For the correctness of our construction, recall that only bits specified in  $\rho_P$  are ever queried when

**P1:** calculating  $g'_{o^*}(x)$ ;      **P2:** verifying  $(f'(x), o^*) \in \mathcal{Q}$ ;      **P3:** verifying  $(x, g'_{o^*}(x)) \notin \mathcal{P}$ .

By our construction, process **S1**, **S2** are sub-procedures of **P1**, **P2**, and thus they will only query locations of  $x$  that are already specified in  $\rho_P$ . However, process **S3** might query some locations that are not specified in  $\rho_P$ . In this case, it is safe to return arbitrary values for those queries. This is because the correctness of reduction  $(h_i, l_o)$  guarantees that there must be  $(y, g_{o^*}(y)) \notin \mathcal{S}$ .

Finally, note that our whole reduction, including the construction of  $(f'_i, g'_o)$  and the execution of process **S1**, **S2**, **S3**, can be done in a uniform manner.  $\square$

With a similar argument, we can also formalize the other direction.

**Lemma 7.7.** *If  $\mathcal{S} \leq_m^U \mathcal{Q}$  and let  $d_1(N), \log M_1(N) = \text{polylog}(N)$ , then*

$$\text{REFUTER}_{d_1, M_1}(\mathcal{P} \rightarrow \mathcal{S}) \leq_m^U \text{REFUTER}_{d_2, M_2}(\mathcal{P} \rightarrow \mathcal{Q}),$$

for some  $d_2(N), \log M_2(N) = \text{polylog}(N)$ .

We often consider all low-depth reductions between two  $\text{TFNP}^{\text{dt}}$  classes with no valid low-depth reductions possible. So, it is convenient to introduce a new kind of  $\text{TFNP}^{\text{dt}}$  subclasses for this type of problem.

**Definition 7.8.** For two  $\text{TFNP}^{\text{dt}}$  classes  $A, B$  ( $A \not\subseteq B$ ) with  $\mathcal{P}, \mathcal{Q}$  being any complete problems of  $A$  and  $B$  respectively,  $\text{Ref}(A \subseteq B)$  is defined as the class of  $\text{TFNP}^{\text{dt}}$  problems that are reducible to  $\text{REFUTER}_{d, M}(\mathcal{P} \rightarrow \mathcal{Q})$  for some  $d(N), \log M(N) = \text{polylog}(N)$ .

This notation is well-defined because [Lemma 7.6](#) and [Lemma 7.7](#) guarantee that the choice of the complete problems does not matter. We also have the following corollary of [Lemma 7.4](#) and [Lemma 7.5](#).

**Corollary 7.9.** *For any two  $\text{TFNP}^{\text{dt}}$  classes  $A, B$  such that  $A \not\subseteq B$ ,  $B \subseteq \text{Ref}(A \subseteq B)$ .*

## 7.2.2 Refuter for Separating from PLS

Now we study the complexity of refuting separations between PLS and other classes in  $\text{TFNP}^{\text{dt}}$ , in particular the separations

$$\text{PPA}^{\text{dt}} \not\subseteq \text{PLS}^{\text{dt}} \quad \text{and} \quad \text{PPP}^{\text{dt}} \not\subseteq \text{PLS}^{\text{dt}}.$$

Our main tool is the equivalence between resolution and PLS via the *false clause search* problem (cf. [\[dRGR22\]](#)): recall that  $\text{Search}(\mathcal{F}) \in \text{PLS}$  if and only if  $\mathcal{F}$  have a  $\text{polylog}(N)$ -width resolution refutation ([Theorem 3.4](#)).

Studying this equivalence from a computational perspective, we related the  $\text{TFNP}^{\text{dt}}$  refuter for PLS with the resolution width refuter.

**Lemma 7.10.** *For any family of unsatisfiable CNF  $\mathcal{F}$  that has no  $\text{polylog}(N)$ -width resolution refutation,*

$$\text{REFUTER}_{d, M}(\text{Search}(\mathcal{F}) \rightarrow \text{ITER}) \leq_m \text{REFUTER}(w(\mathcal{F} \vdash_{\text{Res}} \perp) < w_0)$$

for some  $w_0 = \text{polylog}(N)$  that may depend on  $d, M$ .

Furthermore, this reduction is uniform when  $\mathcal{F}$  is a uniform family of unsatisfiable CNFs.

The proof of [Lemma 7.10](#) follows from the standard procedure that transforms a low-depth decision tree reduction to PLS (i.e., a PLS formulation) to a low-width resolution proof, using the *Prover-Delayer game* [[Pud00](#)]. This proof is rather straightforward, but many details have to be taken care of to make sure that the reduction is uniform. To be self-contained, we formally present the transformation from a PLS formulation to a resolution proof in [Section C.1](#),<sup>33</sup> we then prove [Lemma 7.10](#) in [Section C.2](#).

As an application, we combine [Lemma 7.10](#) with our results on resolution width refuter for EPHP and Tseitin formulas. Note that  $\text{Search}(\text{EPHP})$  and  $\text{Search}(\text{Tseitin})$  are in PPP and PPA respectively. Therefore, we can reduce the  $\text{TFNP}^{\text{dt}}$  refuter for  $\text{PPP}^{\text{dt}} \not\subseteq \text{PLS}^{\text{dt}}$  and  $\text{PPA}^{\text{dt}} \not\subseteq \text{PLS}^{\text{dt}}$  to the resolution width refuters for EPHP and Tseitin respectively.

**Theorem 7.11.** *Let  $\mathcal{P}, \mathcal{Q}$  be any complete problems for PPP and PPA respectively, then for any  $d, M$ , both  $\text{REFUTER}_{d,M}(\mathcal{P} \rightarrow \text{ITER})$  and  $\text{REFUTER}_{d,M}(\mathcal{Q} \rightarrow \text{ITER})$  are PLS-complete via uniform reductions.*

*In particular,  $\text{Ref}(\text{PPP} \subseteq \text{PLS}) = \text{Ref}(\text{PPA} \subseteq \text{PLS}) = \text{PLS}$ .*

Equivalently, [Theorem 7.11](#) says that *local search* arguments are both *necessary* and *sufficient* for separating PPP and PPA from PLS in the black-box setting.

*Proof.* Note that [Lemma 7.5](#) already gives the PLS-hardness result, we will focus on showing that for any  $d, M$ ,  $\text{REFUTER}_{d,M}(\mathcal{P} \rightarrow \text{ITER})$  and  $\text{REFUTER}_{d,M}(\mathcal{Q} \rightarrow \text{ITER})$  are in PLS via uniform reductions.

We start with PPP versus PLS. Note that  $\text{Search}(\text{EPHP})$  is in PPP, thus, by [Lemma 7.6](#), we have

$$\text{REFUTER}_{d,M}(\mathcal{P} \rightarrow \text{ITER}) \leq_m^U \text{REFUTER}_{d',M'}(\text{Search}(\text{EPHP}) \rightarrow \text{ITER}),$$

where  $d', \log M'$  are also poly-logarithmic in  $N$ . Combining with [Lemma 7.10](#), there is

$$\text{REFUTER}_{d,M}(\mathcal{P} \rightarrow \text{ITER}) \leq_m^U \text{REFUTER}(w(\text{EPHP} \vdash_{\text{Res}} \perp) < w_0)$$

for some  $w_0 = \text{polylog}(N)$  that may depend on  $d, M$ .

Recall that [Theorem 4.3](#) shows that  $\text{REFUTER}(w(\text{EPHP} \vdash_{\text{Res}} \perp) < w_0)$  is in PLS via a uniform reduction when  $w_0 = n/3$ . The same reduction to ITER would still work when  $w_0 = \text{polylog}(N)$ , and the only issue is to make sure that the  $\text{cri}(C)$  function could be calculated “efficiently” in this different parameter regime. Note that when  $w_0 = \text{poly}(n)$  (and  $N = 2^{\Omega(n)}$ ), a  $\text{poly}(n)$  time procedure ([Lemma 4.4](#)) would be considered as time efficient; however, when  $w_0 = \text{polylog}(n)$  and  $N$  being quasi-polynomial in  $n$ , only a  $\text{polylog}(n)$  running time is acceptable. Since  $|C| \leq w_0 = \text{polylog}(n)$ , it suffices to prove that the following claim.

**Claim 7.12.**  *$\text{cri}(C)$  can be calculated in  $\text{polylog}(n)$  time when  $|C| = \text{polylog}(n)$ .*

*Proof.* We modify the algorithm described in the proof of [Lemma 4.4](#). First, notice that we do not have to enumerate all possible  $\ell \in [n+1]$ , because only  $\text{polylog}(n)$  pigeons are *involved* in the clause  $C$ , where we say a pigeon  $\ell$  is *involved* in  $C$  if a literal related to  $\ell$  appears in  $C$ . Any pigeons that are not involved in  $C$  would be equivalent, thus, we only need to consider any one of them.

For a fixed  $\ell$ , deciding whether an  $\ell$ -critical assignment exists for  $C$  is reduced to the following graph problem: Given a complete bipartite graph with  $n$  pigeons on the left and  $n$  holes on the right,  $\text{polylog}(n)$  sets of edges are then deleted, determine whether a perfect matching still exists in the end. Each deleted set can be described by a triple  $(i, j_1, j_2)$ , representing the set  $\{(i, j) : j_1 \leq j \leq j_2\}$ .

It is not difficult to design an  $\text{polylog}(n)$  time algorithm for this problem by exploiting the sparsity:

1. We first ignore all pigeons with full degree  $n$ , because they could always be matched in the end.
2. Suppose we have  $t_1 = \text{polylog}(n)$  pigeons left after the first step. We then ignore all pigeons with the degree at least  $t+1$  for the same reason.

<sup>33</sup>This transformation is a well-known folklore among the *Proof Complexity and TFNP* community. However, to the best of the authors’ knowledge, it has not yet been formally written down in any previous literature.

3. We have  $t_2 = \text{polylog}(n)$  pigeons left now, and there are at most  $t_1 \cdot t_2 = \text{polylog}(n)$  edges connected to those pigeons. So, we can run the standard maximum matching algorithm on the subgraph of the remaining pigeons. The original graph has a perfect matching if and only if all  $t_2$  pigeons could be matched.  $\diamond$

A similar argument works for PPA. We use the fact that  $\text{Search}(\text{Tseitin}(G, \tau))$  is in PPA when the graph  $G$  has a constant degree. We will fix a family of strongly explicit expander graph  $G$  and an odd-weighted function  $\tau$ , rather than giving them as input as we did in [Section 6.3](#). For example, we can take  $G$  as a 2D-grid with a boundary being wrapping around, and  $\tau(v) = 1$  only if  $v$  is some designated vertex (say  $(1, 1)$ ). Then, we claim that the  $\text{cri}(C)$  function (defined differently for the Tseitin formula in [Theorem 6.9](#)) can also be calculated in  $\text{polylog}(n)$  time when  $|C| = \text{polylog}(n)$  by exploiting the sparsity of  $C$ . We omit the proof this claim here. Finally, using the same proof of [Theorem 6.9](#), we show that  $\text{REFUTER}(w(\text{Tseitin} \vdash_{\text{Res}} \perp) < w_0$  is in PLS via a uniform reduction when  $w_0 = \text{polylog}(N)$ , which concludes the proof.  $\square$

## Acknowledgments

Jiawei thanks Igor C. Oliveira for introducing him to refuters and their connections with TFNP and thanks Robert Robere and Noah Fleming for knowledge of proof complexity.

Yuhao thanks Toniann Pitassi for the knowledge and guidance in the field of proof complexity and thanks Robert Robere for introducing him to beneficial intuition about proof systems and TFNP<sup>dt</sup> classes.

Hanlin thanks Svyatoslav Gryaznov and Iddo Tzameret for helpful discussions regarding [\[dRGN<sup>+</sup>21\]](#) and proof complexity in general, and Rahul Santhanam and Ján Pich for beneficial conversations.

We thank Lijie Chen, Jiayu Li, and Igor C. Oliveira for sending us a preliminary version of [\[CLO24\]](#). We thank Michal Garlík for helpful discussions on [\[Gar19\]](#). We thank Ján Pich and anonymous referees for their helpful suggestions that improve the presentation of this paper.

## References

- [AB87] Noga Alon and Ravi B. Boppana. The monotone circuit complexity of Boolean functions. *Comb.*, 7(1):1–22, 1987. [doi:10.1007/BF02579196](#). 50
- [ABM23] Albert Atserias, Sam Buss, and Moritz Müller. On the consistency of circuit lower bounds for non-deterministic time. In *STOC*, pages 1257–1270. ACM, 2023. [doi:10.1145/3564246.3585253](#). 13
- [AKPS24] Noel Arteché, Erfan Khaniki, Ján Pich, and Rahul Santhanam. From proof complexity to circuit complexity via interactive protocols. In *ICALP*, volume 297 of *LIPICs*, pages 12:1–12:20. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2024. [doi:10.4230/LIPICs.ICALP.2024.12](#). 1
- [ALWZ21] Ryan Alweiss, Shachar Lovett, Kewen Wu, and Jiapeng Zhang. Improved bounds for the sunflower lemma. *Annals of Mathematics*, 194(3):795–815, 2021. [doi:10.4007/annals.2021.194.3.5](#). 50
- [AM20] Albert Atserias and Moritz Müller. Automating resolution is NP-hard. *J. ACM*, 67(5):31:1–31:17, 2020. [doi:10.1145/3409472](#). 10, 12, 16, 17, 19, 51
- [AT14] Albert Atserias and Neil Thapen. The ordering principle in a fragment of approximate counting. *ACM Trans. Comput. Log.*, 15(4):29:1–29:11, 2014. [doi:10.1145/2629555](#). 7, 24
- [Ats03] Albert Atserias. Improved bounds on the weak pigeonhole principle and infinitely many primes from weaker axioms. *Theor. Comput. Sci.*, 295:27–39, 2003. [doi:10.1016/S0304-3975\(02\)00394-8](#). 6
- [BB17] Arnold Beckmann and Sam Buss. The NP search problems of Frege and Extended Frege proofs. *ACM Trans. Comput. Log.*, 18(2):11:1–11:19, 2017. [doi:10.1145/3060145](#). 14, 66
- [BCE<sup>+</sup>98] Paul Beame, Stephen A. Cook, Jeff Edmonds, Russell Impagliazzo, and Toniann Pitassi. The relative complexity of NP search problems. *J. Comput. Syst. Sci.*, 57(1):3–19, 1998. [doi:10.1006/JCSS.1998.1575](#). 4, 5, 54

- [Bel20] Zoë Bell. Automating regular or ordered resolution is NP-hard. *Electron. Colloquium Comput. Complex.*, TR20-105, 2020. URL: <https://eccc.weizmann.ac.il/report/2020/105>. 12
- [BFI23] Sam Buss, Noah Fleming, and Russell Impagliazzo. TFNP characterizations of proof systems and monotone circuits. In *ITCS*, volume 251 of *LIPICs*, pages 30:1–30:40. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2023. doi:10.4230/LIPICs.ITCS.2023.30. 12, 19, 20, 67
- [BIK<sup>+</sup>94] Paul Beame, Russell Impagliazzo, Jan Krajíček, Toniann Pitassi, and Pavel Pudlák. Lower bound on Hilbert’s Nullstellensatz and propositional proofs. In *FOCS*, pages 794–806. IEEE Computer Society, 1994. doi:10.1109/SFCS.1994.365714. 13
- [BIK<sup>+</sup>97] Samuel R. Buss, Russell Impagliazzo, Jan Krajíček, Pavel Pudlák, Alexander A. Razborov, and Jirí Sgall. Proof complexity in algebraic systems and bounded depth Frege systems with modular counting. *Comput. Complex.*, 6(3):256–298, 1997. doi:10.1007/BF01294258. 13
- [BJ12] Samuel R. Buss and Alan S. Johnson. Propositional proofs and reductions between NP search problems. *Ann. Pure Appl. Log.*, 163(9):1163–1182, 2012. doi:10.1016/J.APAL.2012.01.015. 24, 63
- [BK94] Samuel R. Buss and Jan Krajíček. An application of Boolean complexity to separation problems in bounded arithmetic. *Proceedings of the London Mathematical Society*, s3-69(1):1–21, 1994. doi:10.1112/plms/s3-69.1.1. 4, 7, 8, 10, 14, 24, 25
- [BKO20] Jan Bydzovsky, Jan Krajíček, and Igor C. Oliveira. Consistency of circuit lower bounds with bounded theories. *Log. Methods Comput. Sci.*, 16(2), 2020. doi:10.23638/LMCS-16(2:12)2020. 13
- [BKT14] Samuel R. Buss, Leszek Aleksander Kołodziejczyk, and Neil Thapen. Fragments of approximate counting. *J. Symb. Log.*, 79(2):496–525, 2014. doi:10.1017/JSL.2013.37. 7, 11, 50, 52
- [BKZ15] Samuel R. Buss, Leszek Aleksander Kołodziejczyk, and Konrad Zdanowski. Collapsing modular counting in bounded arithmetic and constant depth propositional proofs. *Trans. Amer. Math. Soc.*, 367(11):7517–7563, 2015. doi:10.1090/S0002-9947-2015-06233-3. 13
- [BM04] Josh Buresh-Oppenheim and Tsuyoshi Morioka. Relativized NP search problems and propositional proof systems. In *CCC*, pages 54–67. IEEE Computer Society, 2004. doi:10.1109/CCC.2004.1313795. 7, 19
- [BM20] Jan Bydzovsky and Moritz Müller. Polynomial time ultrapowers and the consistency of circuit lower bounds. *Arch. Math. Log.*, 59(1-2):127–147, 2020. doi:10.1007/S00153-019-00681-Y. 13
- [BP96] Paul Beame and Toniann Pitassi. Simplified and improved resolution lower bounds. In *FOCS*, pages 274–282. IEEE, 1996. doi:10.1109/SFCS.1996.548486. 1, 3, 8, 9, 14, 15, 18, 26, 28, 29, 39
- [Bro11] L. E. J. Brouwer. Über abbildung von mannigfaltigkeiten. *Mathematische Annalen*, 71:97–115, 1911. In German. doi:10.1007/BF01456931. 4
- [Bus85] Samuel R. Buss. *Bounded arithmetic*. Princeton University, 1985. 4, 14
- [BW01] Eli Ben-Sasson and Avi Wigderson. Short proofs are narrow - resolution made simple. *J. ACM*, 48(2):149–169, 2001. doi:10.1145/375827.375835. 1, 2, 3, 8, 18, 26, 42, 43, 45, 50
- [CDT09] Xi Chen, Xiaotie Deng, and Shang-Hua Teng. Settling the complexity of computing two-player Nash equilibria. *J. ACM*, 56(3):1–57, 2009. doi:10.1145/1516512.1516516. 4
- [CEI96] Matthew Clegg, Jeffery Edmonds, and Russell Impagliazzo. Using the Groebner basis algorithm to find proofs of unsatisfiability. In *STOC*, pages 174–183, 1996. doi:10.1145/237814.237860. 13
- [CJSW24] Lijie Chen, Ce Jin, Rahul Santhanam, and Ryan Williams. Constructive separations and their consequences. *TheoretCS*, volume 3, February 2024. doi:10.46298/theoretics.24.3. 1, 4, 12, 13
- [CK07] Stephen A. Cook and Jan Krajíček. Consequences of the provability of  $\text{NP} \subseteq \text{P}/_{\text{poly}}$ . *J. Symb. Log.*, 72(4):1353–1371, 2007. doi:10.2178/JSL/1203350791. 13
- [CKKO21] Marco Carmosino, Valentine Kabanets, Antonina Kolokolova, and Igor C. Oliveira. LEARN-uniform circuit lower bounds and provability in bounded arithmetic. In *FOCS*, pages 770–780. IEEE, 2021. doi:10.1109/FOCS52979.2021.00080. 13

- [CLO24] Lijie Chen, Jiayu Li, and Igor C. Oliveira. Reverse mathematics of complexity lower bounds. In *FOCS*, pages 505–527. IEEE, 2024. doi:10.1109/FOCS61266.2024.00040. 1, 3, 4, 6, 12, 13, 21, 24, 57
- [CLO25] Lijie Chen, Jiayu Li, and Igor Carboni Oliveira. On the unprovability of circuit size bounds in intuitionistic  $S_2^1$ . *Logical Methods in Computer Science*, Volume 21, Issue 3, Sep 2025. doi:10.46298/lmcs-21(3:26)2025. 3
- [CN10] Stephen A. Cook and Phuong Nguyen. *Logical Foundations of Proof Complexity*, volume 11. Cambridge University Press, 2010. doi:10.1017/CB09780511676277. 3, 21
- [Cob64] Alan Cobham. The intrinsic computational difficulty of functions. In *Proc. Logic, Methodology, and the Philosophy of Science*, pages 24–30, 1964. 21
- [Coo75] Stephen A. Cook. Feasibly constructive proofs and the propositional calculus (preliminary version). In *STOC*, pages 83–97. ACM, 1975. doi:10.1145/800116.803756. 21
- [Coo07] Stephen A. Cook. Bounded reverse mathematics, 2007. Plenary lecture for CiE 2007. 3
- [CP90] Stephen A. Cook and Toniann Pitassi. A feasibly constructive lower bound for resolution proofs. *Inf. Process. Lett.*, 34(2):81–85, 1990. doi:10.1016/0020-0190(90)90141-J. 2, 3, 4, 9, 14, 18
- [CS88] Vasek Chvátal and Endre Szemerédi. Many hard examples for resolution. *J. ACM*, 35(4):759–768, 1988. doi:10.1145/48014.48016. 1, 9, 39, 46, 47, 48, 50, 52
- [CTW23] Lijie Chen, Roei Tell, and Ryan Williams. Derandomization vs refutation: A unified framework for characterizing derandomization. In *FOCS*, pages 1008–1047. IEEE, 2023. doi:10.1109/FOCS57990.2023.00062. 1, 4, 13
- [DGP09] Constantinos Daskalakis, Paul W. Goldberg, and Christos H. Papadimitriou. The complexity of computing a Nash equilibrium. *SIAM J. Comput.*, 39(1):195–259, 2009. doi:10.1137/070699652. 4
- [DLL62] Martin Davis, George Logemann, and Donald W. Loveland. A machine program for theorem-proving. *Commun. ACM*, 5(7):394–397, 1962. doi:10.1145/368273.368557. 1
- [DP60] Martin Davis and Hilary Putnam. A computing procedure for quantification theory. *J. ACM*, 7(3):201–215, 1960. doi:10.1145/321033.321034. 1
- [DR03] Stefan S. Dantchev and Søren Riis. On relativisation and complexity gap for resolution-based proof systems. In *CSL*, volume 2803 of *Lecture Notes in Computer Science*, pages 142–154. Springer, 2003. doi:10.1007/978-3-540-45220-1\_14. 9, 39, 40
- [dRGN<sup>+</sup>21] Susanna F. de Rezende, Mika Göös, Jakob Nordström, Toniann Pitassi, Robert Robere, and Dmitry Sokolov. Automating algebraic proof systems is NP-hard. In *STOC*, pages 209–222. ACM, 2021. doi:10.1145/3406325.3451080. 10, 12, 16, 17, 19, 22, 24, 51, 57
- [dRGR22] Susanna F. de Rezende, Mika Göös, and Robert Robere. Proofs, circuits, and communication. *SIGACT News*, 53(1):59–82, 2022. doi:10.1145/3532737.3532746. 10, 19, 55
- [dRMN<sup>+</sup>20] Susanna F. de Rezende, Or Meir, Jakob Nordström, Toniann Pitassi, Robert Robere, and Marc Vinyals. Lifting with simple gadgets and applications to circuit and proof complexity. In *FOCS*, pages 24–30. IEEE, 2020. doi:10.1109/FOCS46700.2020.00011. 40
- [dRNV16] Susanna F. de Rezende, Jakob Nordström, and Marc Vinyals. How limited interaction hinders real communication (and what it means for proof and circuit complexity). In *FOCS*, pages 295–304. IEEE, 2016. doi:10.1109/FOCS.2016.40. 40
- [Ebt23] Mohammad Hossein Ebtehaj. Variants of pseudo-deterministic algorithms and duality in TFNP. Master’s thesis, University of Waterloo, 2023. URL: <http://hdl.handle.net/10012/19721>. 13
- [ER60] Paul Erdős and Richard Rado. Intersection theorems for systems of sets. *Journal of the London Mathematical Society*, 1(1):85–90, 1960. doi:10.1112/jlms/s1-35.1.85. 50
- [FGPR24] Noah Fleming, Stefan Grosser, Toniann Pitassi, and Robert Robere. Black-box PPP is not Turing-closed. In *STOC*, pages 1405–1414. ACM, 2024. doi:10.1145/3618260.3649769. 24

- [Gar19] Michal Garlík. Resolution lower bounds for refutation statements. In *MFCS*, volume 138 of *LIPICs*, pages 37:1–37:13. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2019. doi:10.4230/LIPICs.MFCS.2019.37.11, 51, 57
- [Gar24] Michal Garlík. Failure of feasible disjunction property for  $k$ -DNF resolution and NP-hardness of automating it. In *CCC*, volume 300 of *LIPICs*, pages 33:1–33:23. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2024. doi:10.4230/LIPICs.CCC.2024.33.12
- [GGKS20] Ankit Garg, Mika Göös, Pritish Kamath, and Dmitry Sokolov. Monotone circuit lower bounds from resolution. *Theory Comput.*, 16:1–30, 2020. doi:10.4086/TOC.2020.V016A013. 40, 50
- [GHJ<sup>+</sup>22] Mika Göös, Alexandros Hollender, Siddhartha Jain, Gilbert Maystre, William Pires, Robert Robere, and Ran Tao. Separations in proof complexity and TFNP. In *FOCS*, pages 1150–1161. IEEE, 2022. doi:10.1109/FOCS54457.2022.00111. 7
- [GKMP20] Mika Göös, Sajin Korothe, Ian Mertz, and Toniann Pitassi. Automating cutting planes is NP-hard. In *STOC*, pages 68–77. ACM, 2020. doi:10.1145/3357713.3384248. 12, 19
- [GP18a] Paul W. Goldberg and Christos H. Papadimitriou. Towards a unified complexity theory of total functions. *J. Comput. Syst. Sci.*, 94:167–192, 2018. doi:10.1016/J.JCSS.2017.12.003. 12, 14, 66
- [GP18b] Mika Göös and Toniann Pitassi. Communication lower bounds via critical block sensitivity. *SIAM J. Comput.*, 47(5):1778–1806, 2018. doi:10.1137/16M1082007. 40
- [GST07] Dan Gutfreund, Ronen Shaltiel, and Amnon Ta-Shma. If NP languages are hard on the worst-case, then it is easy to find their hard instances. *Comput. Complex.*, 16(4):412–441, 2007. doi:10.1007/S00037-007-0235-8. 4
- [Hak85] Armin Haken. The intractability of resolution. *Theor. Comput. Sci.*, 39:297–308, 1985. doi:10.1016/0304-3975(85)90144-6. 1, 2, 3, 5, 9, 18, 28, 29
- [Hak95] Armin Haken. Counting bottlenecks to show monotone  $P \neq NP$ . In *FOCS*, pages 36–40. IEEE Computer Society, 1995. doi:10.1109/SFCS.1995.492460. 50
- [HN12] Trinh Huynh and Jakob Nordström. On the virtue of succinct proofs: Amplifying communication complexity hardness to time-space trade-offs in proof complexity. In *STOC*, pages 233–248. ACM, 2012. doi:10.1145/2213977.2214000. 40
- [ILW23] Rahul Ilango, Jiayu Li, and R. Ryan Williams. Indistinguishability obfuscation, range avoidance, and bounded arithmetic. In *STOC*, pages 1076–1089. ACM, 2023. doi:10.1145/3564246.3585187. 6
- [IR22] Dmitry Itsykson and Artur Riazanov. Automating OBDD proofs is NP-hard. In *MFCS*, volume 241 of *LIPICs*, pages 59:1–59:15. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2022. doi:10.4230/LIPICs.MFCS.2022.59.12
- [Jeř04] Emil Jeřábek. Dual weak pigeonhole principle, Boolean complexity, and derandomization. *Ann. Pure Appl. Log.*, 129(1-3):1–37, 2004. doi:10.1016/j.apal.2003.12.003. 6, 21, 24, 25
- [Jeř07a] Emil Jeřábek. Approximate counting in bounded arithmetic. *J. Symb. Log.*, 72(3):959–993, 2007. doi:10.2178/JSL/1191333850. 6, 50
- [Jeř07b] Emil Jeřábek. On independence of variants of the weak pigeonhole principle. *J. Log. Comput.*, 17(3):587–604, 2007. doi:10.1093/LOGCOM/EXM017. 6, 24
- [Jeř09] Emil Jeřábek. Approximate counting by hashing in bounded arithmetic. *J. Symb. Log.*, 74(3):829–860, 2009. doi:10.2178/JSL/1245158087. 7, 50
- [JPY88] David S. Johnson, Christos H. Papadimitriou, and Mihalis Yannakakis. How easy is local search? *J. Comput. Syst. Sci.*, 37(1):79–100, 1988. doi:10.1016/0022-0000(88)90046-3. 4, 8, 19
- [Kam19] Pritish Kamath. *Some hardness escalation results in computational complexity theory*. PhD thesis, Massachusetts Institute of Technology, 2019. URL: <https://hdl.handle.net/1721.1/128290>. 10
- [KKMP21] Robert Kleinberg, Oliver Korten, Daniel Mitropolsky, and Christos H. Papadimitriou. Total functions in the polynomial hierarchy. In *ITCS*, volume 185 of *LIPICs*, pages 44:1–44:18, 2021. doi:10.4230/LIPICs.ITCS.2021.44.6

- [KNT11] Leszek Aleksander Kołodziejczyk, Phuong Nguyen, and Neil Thapen. The provably total NP search problems of weak second order bounded arithmetic. *Ann. Pure Appl. Log.*, 162(6):419–446, 2011. doi:10.1016/J.APAL.2010.12.002. 14
- [KO17] Jan Krajíček and Igor C. Oliveira. Unprovability of circuit upper bounds in Cook’s theory PV. *Log. Methods Comput. Sci.*, 13(1), 2017. doi:10.23638/LMCS-13(1:4)2017. 13
- [Kor21] Oliver Korten. The hardest explicit construction. In *FOCS*, pages 433–444. IEEE, 2021. doi:10.1109/FOCS52979.2021.00051. 6, 24
- [Kor22] Oliver Korten. Derandomization from time-space tradeoffs. In *CCC*, volume 234 of *LIPICs*, pages 37:1–37:26. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2022. doi:10.4230/LIPICs.CCC.2022.37.1, 4, 6, 12, 13, 24
- [KP89] Jan Krajíček and Pavel Pudlák. Propositional provability and models of weak arithmetic. In *CSL*, volume 440 of *Lecture Notes in Computer Science*, pages 193–210. Springer, 1989. doi:10.1007/3-540-52753-2\_40. 3
- [Kra95] Jan Krajíček. *Bounded arithmetic, propositional logic, and complexity theory*, volume 60 of *Encyclopedia of mathematics and its applications*. Cambridge University Press, 1995. doi:10.1017/CB09780511529948. 20, 21
- [Kra97] Jan Krajíček. Interpolation theorems, lower bounds for proof systems, and independence results for bounded arithmetic. *J. Symb. Log.*, 62(2):457–486, 1997. doi:10.2307/2275541. 3, 50
- [Kra01] Jan Krajíček. On the weak pigeonhole principle. *Fundamenta Mathematicae*, 170(1-2):123–140, 2001. URL: <http://eudml.org/doc/282141>. 6
- [Kra04] Jan Krajíček. Dual weak pigeonhole principle, pseudo-surjective functions, and provability of circuit lower bounds. *J. Symb. Log.*, 69(1):265–286, 2004. doi:10.2178/jsl/1080938841. 6, 24
- [Kra11a] Jan Krajíček. A note on propositional proof complexity of some Ramsey-type statements. *Arch. Math. Log.*, 50(1-2):245–255, 2011. doi:10.1007/S00153-010-0212-9. 39, 40
- [Kra11b] Jan Krajíček. On the proof complexity of the Nisan-Wigderson generator based on a hard  $\text{NP} \cap \text{coNP}$  function. *J. Math. Log.*, 11(1), 2011. doi:10.1142/S0219061311000979. 3
- [Kra19] Jan Krajíček. *Proof Complexity*. Encyclopedia of Mathematics and its Applications. Cambridge University Press, 2019. doi:10.1017/9781108242066. 1, 5, 40
- [KST07] Jan Krajíček, Alan Skelley, and Neil Thapen. NP search problems in low fragments of bounded arithmetic. *J. Symb. Log.*, 72(2):649–672, 2007. doi:10.2178/JSL/1185803628. 14
- [KT22] Leszek Aleksander Kołodziejczyk and Neil Thapen. Approximate counting and NP search problems. *J. Math. Log.*, 22(3):2250012:1–2250012:31, 2022. doi:10.1142/S021906132250012X. 14
- [LC11] Dai Tri Man Le and Stephen A. Cook. Formalizing randomized matching algorithms. *Log. Methods Comput. Sci.*, 8(3), 2011. doi:10.2168/LMCS-8(3:5)2012. 27
- [LO23] Jiayu Li and Igor C. Oliveira. Unprovability of strong complexity lower bounds in bounded arithmetic. In *STOC*, pages 1051–1057. ACM, 2023. doi:10.1145/3564246.3585144. 3
- [Maa84] Wolfgang Maass. Quadratic lower bounds for deterministic and nondeterministic one-tape Turing machines (extended abstract). In *STOC*, pages 401–408. ACM, 1984. doi:10.1145/800057.808706. 3
- [Mor01] Tsuyoshi Morioka. Classification of search problems and their definability in bounded arithmetic. Master’s thesis, University of Toronto, 2001. URL: <https://hdl.handle.net/1807/16458>. 19
- [MP91] Nimrod Megiddo and Christos H. Papadimitriou. On total functions, existence theorems and computational complexity. *Theor. Comput. Sci.*, 81(2):317–324, 1991. doi:10.1016/0304-3975(91)90200-L. 1, 4
- [MP96] Alexis Maciel and Toniann Pitassi. Towards lower bounds for bounded-depth Frege proofs with modular connectives. In *Proof Complexity and Feasible Arithmetics*, volume 39 of *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, pages 195–227. DIMACS/AMS, 1996. doi:10.1090/DIMACS/039/12. 13

- [MP20] Moritz Müller and Ján Pich. Feasibly constructive proofs of succinct weak circuit lower bounds. *Ann. Pure Appl. Log.*, 171(2), 2020. doi:10.1016/j.apal.2019.102735. 3, 4
- [MPW02] Alexis Maciel, Toniann Pitassi, and Alan R. Woods. A new proof of the weak pigeonhole principle. *J. Comput. Syst. Sci.*, 64(4):843–872, 2002. doi:10.1006/JCSS.2002.1830. 6
- [Mül21] Moritz Müller. Typical forcings, NP search problems and an extension of a theorem of Riis. *Ann. Pure Appl. Log.*, 172(4):102930, 2021. doi:10.1016/J.APAL.2020.102930. 4, 27
- [Nas51] John Nash. Non-cooperative games. *Annals of mathematics*, 54(2):286–295, 1951. doi:10.2307/1969529. 4
- [Ngu08] Phuong Nguyen. *Bounded reverse mathematics*. PhD thesis, University of Toronto, 2008. 3
- [Pap94] Christos H. Papadimitriou. On the complexity of the parity argument and other inefficient proofs of existence. *J. Comput. Syst. Sci.*, 48(3):498–532, 1994. doi:10.1016/S0022-0000(05)80063-7. 4
- [Pap24] Theodoros Papamakarios. Depth- $d$  frege systems are not automatable unless  $P = NP$ . In *CCC*, volume 300 of *LIPICs*, pages 22:1–22:17. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2024. doi:10.4230/LIPICs.CCC.2024.22. 12
- [Pic15] Ján Pich. Circuit lower bounds in bounded arithmetics. *Ann. Pure Appl. Log.*, 166(1):29–45, 2015. doi:10.1016/J.APAL.2014.08.004. 3, 4
- [PS19] Ján Pich and Rahul Santhanam. Why are proof complexity lower bounds hard? In *FOCS*, pages 1305–1324. IEEE Computer Society, 2019. doi:10.1109/FOCS.2019.00080. 1
- [PS21] Ján Pich and Rahul Santhanam. Strong co-nondeterministic lower bounds for NP cannot be proved feasibly. In *STOC*, pages 223–233. ACM, 2021. doi:10.1145/3406325.3451117. 3
- [PS26] Ján Pich and Rahul Santhanam. Towards  $P \neq NP$  from Extended Frege lower bounds. *J. ACM*, 2026. doi:10.1145/3801091. 1, 13
- [PT12] Pavel Pudlák and Neil Thapen. Alternating minima and maxima, Nash equilibria and Bounded Arithmetic. *Ann. Pure Appl. Log.*, 163(5):604–614, 2012. doi:10.1016/J.APAL.2011.06.014. 14
- [PT19] Pavel Pudlák and Neil Thapen. Random resolution refutations. *Comput. Complex.*, 28(2):185–239, 2019. doi:10.1007/S00037-019-00182-7. 2, 7, 11, 24, 50, 52
- [Pud97] Pavel Pudlák. Lower bounds for resolution and cutting plane proofs and monotone computations. *J. Symb. Log.*, 62(3):981–998, 1997. doi:10.2307/2275583. 50
- [Pud00] Pavel Pudlák. Proofs as games. *Am. Math. Mon.*, 107(6):541–550, 2000. doi:10.2307/2589349. 51, 56, 68, 70
- [Pud20] Pavel Pudlák. Reflection principles, propositional proof systems, and theories. *arXiv*, abs/2007.14835, 2020. doi:10.48550/arXiv.2007.14835. 14, 67
- [PWW88] Jeff B. Paris, A. J. Wilkie, and Alan R. Woods. Provability of the pigeonhole principle and the existence of infinitely many primes. *J. Symb. Log.*, 53(4):1235–1244, 1988. doi:10.1017/S0022481200028061. 6, 21, 24
- [Raz85] Alexander A. Razborov. Lower bounds on the monotone complexity of some Boolean function. In *Soviet Math. Dokl.*, volume 31, pages 354–357, 1985. URL: <https://www.mathnet.ru/eng/dan9192>. 50
- [Raz87] Alexander A. Razborov. Lower bounds on the size of bounded depth circuits over a complete basis with logical addition. *Mathematical Notes of the Academy of Sciences of the USSR*, 41(4):333–338, 1987. 13
- [Raz95a] Alexander A. Razborov. Bounded arithmetic and lower bounds in Boolean complexity. In *Feasible Mathematics II*, pages 344–386. Birkhäuser Boston, 1995. doi:10.1007/978-1-4612-2566-9\_12. 3
- [Raz95b] Alexander A. Razborov. Unprovability of lower bounds on circuit size in certain fragments of bounded arithmetic. *Izvestiya: mathematics*, 59(1):205, 1995. doi:10.1070/IM1995v059n01ABEH000009. 3, 10, 50

- [Raz98] Alexander A. Razborov. Lower bounds for the polynomial calculus. *Comput. Complex.*, 7(4):291–324, 1998. doi:10.1007/S000370050013. 13
- [Raz15] Alexander A. Razborov. Pseudorandom generators hard for  $k$ -DNF resolution and polynomial calculus resolution. *Annals of Mathematics*, 181(2):415–472, 2015. doi:10.4007/annals.2015.181.2.1. 1
- [Sch97] Uwe Schöning. Resolution proofs, exponential bounds, and Kolmogorov complexity. In *MFCS*, volume 1295 of *Lecture Notes in Computer Science*, pages 110–116. Springer, 1997. doi:10.1007/BFB0029954. 1, 9, 15, 39, 42, 45
- [Smo87] Roman Smolensky. Algebraic methods in the theory of lower bounds for boolean circuit complexity. In *STOC*, pages 77–82. ACM, 1987. doi:10.1145/28395.28404. 13
- [ST11] Alan Skelley and Neil Thapen. The provably total search problems of bounded arithmetic. *Proceedings of the London Mathematical Society*, 103(1):106–138, 2011. doi:10.1112/plms/pdq044. 14
- [ST21] Rahul Santhanam and Iddo Tzameret. Iterated lower bound formulas: a diagonalization-based approach to proof complexity. In *STOC*, pages 234–247. ACM, 2021. doi:10.1145/3406325.3451010. 1
- [SW14] Rahul Santhanam and Ryan Williams. On uniformity and circuit lower bounds. *Comput. Complex.*, 23(2):177–205, 2014. doi:10.1007/S00037-014-0087-Y. 13
- [Tha02] Neil Thapen. *The weak pigeonhole principle in models of bounded arithmetic*. PhD thesis, University of Oxford, 2002. 6, 24
- [Tse83] Grigori S Tseitin. On the complexity of derivation in propositional calculus. *Automation of reasoning: 2: Classical papers on computational logic 1967–1970*, pages 466–483, 1983. doi:10.1007/978-3-642-81955-1\_28. 42
- [Urq87] Alasdair Urquhart. Hard examples for resolution. *J. ACM*, 34(1):209–219, 1987. doi:10.1145/7531.8928. 1, 9, 39, 42
- [VS23] Bernhard Von Stengel. Zero-sum games and linear programming duality. *Mathematics of Operations Research*, 2023. doi:10.1287/moor.2022.0149. 4

## A Amplification for $\text{rwPHP}(\mathcal{P})$

In this section, we prove [Theorem 3.12](#), showing that the relationship between  $N$  and  $M$  does not influence the complexity of  $\text{rwPHP}(\mathcal{P})$  (provided that  $M$  and  $N$  are not too close to each other), thus our choice of  $N = 2M$  is indeed without loss of generality. This result requires  $\mathcal{P}$  to be closed under Turing reductions [[BJ12](#)], as defined below:

**Definition A.1** (Turing Reductions in  $\text{TFNP}^{\text{dt}}$ ). Let  $\mathcal{P}, \mathcal{Q}$  be problems in  $\text{TFNP}^{\text{dt}}$ . We say there is a time- $t$  (uniform) Turing reduction from  $\mathcal{Q}$  to  $\mathcal{P}$  if there is a time- $t$  oracle Turing machine  $R^{x, \mathcal{P}}$  that solves  $\mathcal{Q}$  in the following manner. Let  $x \in \{0, 1\}^N$  be the input to  $\mathcal{Q}$ . Besides work tapes and a query tape for access to  $x$ ,  $R$  has another query tapes for access to a  $\mathcal{P}$  oracle. Each query  $q$  to  $\mathcal{P}$  is described as  $(1^{t'}, L, M_q^x)$ , where  $L$  is the length of the query input, and  $M_q^x$  is a time- $t'$  Turing machine with query access to  $x$ . The input to  $\mathcal{P}$  is defined as the  $L$ -bit string whose  $i$ -th bit is  $M_q^x(i)$ . The answer to this query would be any valid solution for this  $L$ -bit string as an input to  $\mathcal{P}$ . Finally, for every  $x \in \{0, 1\}^N$  and every valid computational history of  $R$  (i.e., every query to  $\mathcal{P}$  is answered correctly), the output of  $R$  should be a valid output of  $x$  for  $\mathcal{Q}$ .

**Assumption A.2.**  $\mathcal{P}$  is closed under Turing reductions. More precisely, for some function  $\gamma(t)$  (think of  $\gamma(t) \leq \text{poly}(t)$ ), if a  $\text{TFNP}^{\text{dt}}$  problem  $\mathcal{Q}$  admits a time- $t$  Turing reduction to  $\mathcal{P}$ , then  $\mathcal{Q}$  also admits a uniform depth- $\gamma(t)$  decision tree reduction to  $\mathcal{P}$ .

Recall that  $\text{rPHP}(\mathcal{P})_{M \rightarrow N}$  denotes the  $\text{rwPHP}(\mathcal{P})$  problem where the purported “surjection” is  $f : [M] \rightarrow [N]$ .

**Fact A.3.** *Let  $M < N_1 \leq N_2$ , then there is a depth-1 decision tree reduction from  $\text{rPHP}(\mathcal{P})_{M \rightarrow N_2}$  to  $\text{rPHP}(\mathcal{P})_{M \rightarrow N_1}$ .*

**Theorem A.4** (Formal version of [Theorem 3.12](#)). *Let  $N \geq 2M$  and  $\varepsilon > 0$  be parameters, and let  $d := \gamma(O(\varepsilon^{-1})) \cdot \gamma(O(\log \frac{N}{M}))$ . If [Assumption A.2](#) holds for  $\mathcal{P}$ , then there is a depth- $d$  decision tree reduction from  $\text{rPHP}(\mathcal{P})_{M \rightarrow (1+\varepsilon)M}$  to  $\text{rPHP}(\mathcal{P})_{M \rightarrow N}$ .*

We prove [Theorem A.4](#) in two steps: in [Lemma A.5](#) we reduce  $\text{rPHP}$  with stretch  $(1 + \varepsilon)$  (i.e.,  $\text{rPHP}(\mathcal{P})_{M \rightarrow (1+\varepsilon)M}$ ) to  $\text{rwPHP}$  with stretch 2, and in [Lemma A.6](#) we reduce  $\text{rwPHP}$  with stretch 2 to  $\text{rwPHP}$  with arbitrarily large stretch. [Theorem A.4](#) follows easily from [Lemma A.5](#) and [A.6](#).

**Lemma A.5.** *Let  $M \geq 1$ ,  $\varepsilon > 0$  be parameters. Suppose that [Assumption A.2](#) holds for  $\mathcal{P}$ . Then there is a depth- $\gamma(O(\varepsilon^{-1}))$  decision tree reduction from  $\text{rPHP}(\mathcal{P})_{M \rightarrow \lfloor (1+\varepsilon)M \rfloor}$  to  $\text{rPHP}(\mathcal{P})_{M \rightarrow 2M}$ .*

*Proof.* Without loss of generality, assume that both  $\varepsilon \cdot M$  and  $d := 1/\varepsilon$  are integers. Let  $(f, \{I_y\}, \{g_y\})$  be an instance of  $\text{rPHP}(\mathcal{P})_{M \rightarrow (1+\varepsilon)M}$  and we want to reduce it to an instance  $(f', \{I'_y\}, \{g'_y\})$  of  $\text{rPHP}(\mathcal{P})_{M \rightarrow 2M}$ . Recall that:

- $f : [M] \rightarrow [(1 + \varepsilon)M]$  is the purported “surjection”.
- For every  $y \in [(1 + \varepsilon)M]$ ,  $I_y$  is a  $\mathcal{P}$  instance where every possible answer  $ans$  of  $I_y$  is labelled with an integer  $g_y(ans) \in [M]$ .
- The goal is to find some  $y \in [(1 + \varepsilon)M]$  and a solution  $ans$  of  $I_y$  such that  $f(g_y(ans)) \neq y$ .

For every  $k \in [d]$  (recall  $d = 1/\varepsilon$ ), define  $f_k : [M + k\varepsilon M] \rightarrow [M + (k + 1)\varepsilon M]$  as the following function: on input  $x \in [M + k\varepsilon M]$ , if  $x < M$ , then  $f_k(x) := f(x)$ ; otherwise  $f_k(x) := x + \varepsilon M$ . The function  $f'$  in our reduction is simply  $f' := f_{d-1} \circ f_{d-2} \circ \dots \circ f_0$ . Intuitively, if (a weak theory thinks that)  $f : [M] \rightarrow [(1 + \varepsilon)M]$  is a surjection, then (it also thinks that)  $f' : [M] \rightarrow [2M]$  is a surjection.

Next we define the instances  $\{I'_y\}$  and the functions  $\{g'_y\}$ . Roughly speaking, the input instance  $\{I_y\}$  and  $\{g_y\}$  defines a  $\mathcal{P}$ -computable multi-function (also denoted as)  $g : [(1 + \varepsilon)M] \rightarrow [M]$ , which is a purported inverse of  $f$ . By padding  $g$ , we obtain  $\mathcal{P}$ -computable multi-functions  $g_k : [M + (k + 1)\varepsilon M] \rightarrow [M + k\varepsilon M]$  for each  $k \in [d]$ , and each  $g_k$  is a purported inverse of  $f_k$ . We compose these multi-functions  $g_k$  to obtain a single multi-function  $g : [2M] \rightarrow [M]$  that can be computed by a Turing reduction to  $\mathcal{P}$ . Details follow.

Consider a Turing machine with oracle access to  $I_y$  and  $\mathcal{P}$  that on input  $y \in [2M]$ , operates as follows. Let  $y_d := y$ . For each  $k$  from  $d - 1$  downto 0:

- if  $y_{k+1} \geq (1 + \varepsilon)M$ , then we define  $y_k := y_{k+1} - \varepsilon M$ ;
- otherwise we query  $\mathcal{P}$  to obtain a valid answer  $ans_k$  for  $I_{y_{k+1}}$  and let  $y_k := g(ans_k)$ .

Finally, the machine outputs  $y_0$  (as a purported preimage of  $y$  under  $f'$ ).

The computational history of this Turing machine defines a total search problem  $L_{\text{his}}$  as follows. The input to  $L_{\text{his}}$  consists of  $(f, \{I_y\}, \{g_y\})$ , as well as some  $y \in [2M]$ . The output consists of a sequence  $(ans_0, ans_1, \dots, ans_{d-1})$ . Denote  $y_d = y$  and

$$y_k = \begin{cases} y_{k+1} - \varepsilon M & \text{if } y_{k+1} \geq (1 + \varepsilon)M, \\ g_{y_{k+1}}(ans_k) & \text{otherwise} \end{cases}$$

for each  $k \in [d]$ . We accept the output if for every  $k$  such that  $y_{k+1} < (1 + \varepsilon)M$ ,  $ans_k$  is a valid solution for  $I_{y_{k+1}}$ ; otherwise we reject the output.

Clearly, the above Turing machine itself is a time- $O(d)$  Turing reduction from  $L_{\text{his}}$  to  $\mathcal{P}$ . Since  $\mathcal{P}$  is closed under Turing reductions, there is also a depth- $\gamma(O(d))$  mapping reduction from  $L_{\text{his}}$  to  $\mathcal{P}$ . That

is, there is a depth- $\gamma(O(d))$  decision tree that on input  $(f, \{I_y\}, \{g_y\})$  as well as  $y \in [2M]$ , outputs a  $\mathcal{P}$  instance (that we call  $I'_y$ ), and a mapping that given any valid answer  $ans$  of  $I'_y$ , finds a valid sequence  $(ans_0, ans_1, \dots, ans_{d-1})$ . We compute each  $\{y_k\}_{k \in [d+1]}$  as above and define  $g'_y(ans) := y_0$ .

This finishes the description of our reduction from  $\text{rPHP}(\mathcal{P})_{M \rightarrow (1+\varepsilon)M}$  to  $\text{rPHP}(\mathcal{P})_{M \rightarrow 2M}$ ; it is easy to see that it has depth  $\gamma(O(\varepsilon^{-1}))$ . Now, given a valid solution  $(y', ans')$  for  $(f', \{I'_y\}, \{g'_y\})$ , we can compute a valid solution  $(y, ans)$  for  $(f, \{I_y\}, \{g_y\})$  as follows. First, since  $ans'$  is a valid solution for  $I'_{y'}$ , we can unpack  $ans'$  to obtain a sequence  $(ans_0, ans_1, \dots, ans_{d-1})$ . Then we define each  $\{y_k\}_{k \in [d+1]}$  as above (starting with  $y_d = y'$ ). Also, for every  $k \in [d+1]$ , define  $f_{\geq k} := f_{d-1} \circ \dots \circ f_k$ , then  $f_{\geq k}$  is a purported surjection from  $[M + k\varepsilon M]$  to  $[2M]$ . (As special cases,  $f_{\geq d} : [2M] \rightarrow [2M]$  is the identity function and  $f_{\geq 0} = f'$ .) Since  $(y', ans')$  is a valid solution, we know that  $f'(g'_{y'}(ans')) \neq y'$ , which translates to  $f_{\geq 0}(y_0) \neq y_d$ . Since  $f_{\geq d}(y_d) = y_d$ , there is some integer  $k \in [d]$  such that  $f_{\geq k}(y_k) \neq y_d$  but  $f_{\geq k+1}(y_{k+1}) = y_d$ . We argue that  $(y_{k+1}, ans_k)$  is a valid solution for  $(f, \{I_y\}, \{g_y\})$ :

- First, it must be the case that  $y_{k+1} < (1 + \varepsilon)M$ . If  $y_{k+1} \geq (1 + \varepsilon)M$ , then  $y_k = y_{k+1} - \varepsilon M \geq M$  and thus  $f_k(y_k) = y_{k+1}$ . It follows that

$$y_d \neq f_{\geq k}(y_k) = f_{\geq k+1}(f_k(y_k)) = f_{\geq k+1}(y_{k+1}) = y_d, \quad (9)$$

a contradiction.

- Since  $(ans_0, ans_1, \dots, ans_{d-1})$  is a valid sequence,  $ans_k$  is a valid solution for  $I_{y_{k+1}}$ .
- Finally, if  $f(g(ans_k)) = f(y_k) = y_{k+1}$ , then (9) holds, which is a contradiction. Therefore, it must be the case that  $f(g_{y_{k+1}}(ans_k)) \neq y_{k+1}$  and thus  $(y_{k+1}, ans_k)$  is a valid solution for  $(f, \{I_y\}, \{g_y\})$ .  $\square$

**Lemma A.6.** *Let  $N \geq 2M$ . Suppose that [Assumption A.2](#) holds for  $\mathcal{P}$ . There is a depth- $\gamma(O(\log \frac{N}{M}))$  decision tree reduction from  $\text{rPHP}(\mathcal{P})_{M \rightarrow 2M}$  to  $\text{rPHP}(\mathcal{P})_{M \rightarrow N}$ .*

*Proof.* The proof is similar to that of [Lemma A.5](#). Without loss of generality, we may assume that  $d := \log \frac{N}{M}$  is an integer (i.e.,  $N/M$  is a power of 2). Let  $(f, \{I_y\}, \{g_y\})$  be an instance of  $\text{rPHP}(\mathcal{P})_{M \rightarrow 2M}$  and we want to reduce it to an instance  $(f', \{I'_y\}, \{g'_y\})$  of  $\text{rPHP}(\mathcal{P})_{M \rightarrow N}$ . Recall that:

- $f : [M] \rightarrow [2M]$  is the purported “surjection”.
- For every  $y \in [2M]$ ,  $I_y$  is a  $\mathcal{P}$  instance where every possible answer  $ans$  of  $I_y$  is labelled with an integer  $g(ans) \in [M]$ .
- The goal is to find some  $y \in [2M]$  and a solution  $ans$  of  $I_y$  such that  $f(g_y(ans)) \neq y$ .

For every integer  $k \in [d]$ , we put  $2^k$  copies of the instance  $(f, \{I_y\}, \{g_y\})$  in parallel and obtain the instance  $(f_k, \{(I_k)_y\}, \{g_{k,y}\})$  of  $\text{rPHP}(\mathcal{P})_{(2^k M) \rightarrow (2^{k+1} M)}$ . More precisely:

(Definition of  $f_k$ ) Let  $x \in [2^k M]$  be the input, and let  $x = x_0 \cdot M + x_1$  where  $x_0 \in [2^k]$  and  $x_1 \in [M]$ . We define  $f_k(x) := x_0 \cdot 2M + f(x_1)$ .

(Definition of  $(I_k)_y$  and  $g_{k,y}$ ) Let  $y \in [2^{k+1} M]$  be the input, and let  $y = y_0 \cdot 2M + y_1$  where  $y_0 \in [2^k]$  and  $y_1 \in [2M]$ . We define  $(I_k)_y := I_{y_1}$ , and for every  $ans$  that is a possible solution of  $(I_k)_y = I_{y_1}$ , define  $g_{k,y}(ans) := y_0 \cdot M + g_y(ans)$ .

The mapping from  $(f, \{I_y\}, \{g_y\})$  to  $(f_k, \{(I_k)_y\}, \{g_{k,y}\})$  can be computed by a depth-1 decision tree. Given a valid solution  $(y, ans)$  for  $(f_k, \{(I_k)_y\}, \{g_{k,y}\})$ , we write  $y = y_0 \cdot 2M + y_1$  where  $y_0 \in [2^k]$  and  $y_1 \in [2M]$ . Since  $ans$  is a solution of  $(I_k)_y = I_{y_1}$  and

$$y_0 \cdot 2M + y_1 = y \neq f_k(g_{k,y}(ans)) = y_0 \cdot 2M + f(g(ans)) \implies f(g_y(ans)) \neq y_1,$$

it follows that  $(y_1, ans)$  is also a valid solution for  $(f, \{I_y\}, \{g_y\})$ . Therefore, there is a depth-1 decision tree reduction from  $\text{rPHP}(\mathcal{P})_{M \rightarrow 2M}$  to  $\text{rPHP}(\mathcal{P})_{(2^k M) \rightarrow (2^{k+1} M)}$ .

Now, we compose the instances  $(f_k, \{(I_k)_y\}, g_{k,y})$  for every  $k \in [d]$  to obtain the instance  $(f', \{I'_y\}, \{g'_y\})$ . In particular, the ‘‘surjection’’  $f' : [M] \rightarrow [2^d M]$  is defined as  $f' := f_{d-1} \circ f_{d-2} \circ \dots \circ f_0$ .

To define  $I'_y$  and  $g'_y$ , consider the Turing machine with oracle access to  $\mathcal{P}$  that, on input  $y \in [2^d M]$ , operates as follows. Let  $y_d := y$ . For each  $k$  from  $d-1$  to  $0$ , the machine queries  $\mathcal{P}$  to obtain a valid answer  $ans_k$  for  $(I_k)_{y_{k+1}}$ , and then sets  $y_k := g_{k,y_{k+1}}(ans_k)$ . Finally, the machine outputs the number  $y_0 \in [M]$ .

We define a total search problem  $L_{\text{his}}$  based on the computational history of this machine. The input of  $L_{\text{his}}$  consists of  $M, N, (f, \{I_y\}, \{g_y\})$ , as well as some  $y \in [2^d M]$ ; note that given these inputs, one can define the  $\text{rPHP}(\mathcal{P})_{(2^k M) \rightarrow (2^{k+1} M)}$  instances  $(f_k, \{(I_k)_y\}, \{g_{k,y}\})$  as before. The output consists of a sequence  $(ans_0, ans_1, \dots, ans_{d-1})$ . Denoting  $y_d = y$  and  $y_k = g_{k,y_{k+1}}(ans_k)$  for every  $k \in [d]$ , accept the output if for every  $k \in [d]$ ,  $ans_k$  is a valid solution for  $(I_k)_{y_{k+1}}$ ; otherwise reject the output.

Clearly, the above Turing machine itself is a time- $O(d)$  Turing reduction from  $L_{\text{his}}$  to  $\mathcal{P}$ . Since  $\mathcal{P}$  is closed under Turing reductions, there is also a depth- $\gamma(O(d))$  mapping reduction from  $L_{\text{his}}$  to  $\mathcal{P}$ . Therefore, there is a depth- $\gamma(O(d))$  decision tree that on input  $(f, \{I_y\}, \{g_y\})$  as well as  $y \in [2^d M]$ , outputs a  $\mathcal{P}$  instance (that we call  $I'_y$ ), and a mapping that given any valid answer  $ans$  of  $I'_y$ , finds a valid sequence  $(ans_0, ans_1, \dots, ans_{d-1})$ . We define  $g'_y(ans) := g_{0,y}(ans_0)$ .

This finishes the description of our reduction from  $\text{rPHP}(\mathcal{P})_{M \rightarrow 2M}$  to  $\text{rPHP}(\mathcal{P})_{M \rightarrow N}$ ; it is easy to see that it has depth  $\gamma(O(\log \frac{N}{M}))$ . Now, given a valid solution  $(y', ans')$  for  $(f', \{I'_y\}, \{g'_y\})$ , we can compute a valid solution  $(y, ans)$  for  $(f, \{I_y\}, \{g_y\})$  as follows. First, since  $ans'$  is a valid solution for  $I'_{y'}$ , we can unpack  $ans'$  to obtain a sequence  $ans_0, ans_1, \dots, ans_{d-1}$ . Let  $y_d = y'$  and  $y_k = g_{k,y_{k+1}}(ans_k)$  for every  $k$  from  $d-1$  down to  $0$ , then  $f'(y_0) \neq y_d$ . For every  $k \in \{0, 1, \dots, d\}$ , let  $f_{\geq k} := f_{d-1} \circ f_{d-2} \circ \dots \circ f_k$ ; notice that  $f_{\geq k}$  is a purported surjection from  $[2^k M]$  to  $[2^d M]$ . (Note that as special cases,  $f_{\geq 0} = f'$  and  $f_{\geq d} : [2^d M] \rightarrow [2^d M]$  is the identity function.) Since  $f_{\geq 0}(y_0) \neq y_d$  but  $f_{\geq d}(y_d) = y_d$ , there is an integer  $k \in [d]$  such that  $f_{\geq k}(y_k) \neq y_d$  but  $f_{\geq (k+1)}(y_{k+1}) = y_d$ . We claim that  $(y_{k+1}, ans_k)$  is a valid solution to the instance  $(f_k, \{(I_k)_y\}, \{g_{k,y}\})$ .

- Since  $(ans_0, ans_1, \dots, ans_{d-1})$  is a valid solution of  $L_{\text{his}}$  on input  $(f, \{I_y\}, \{g_y\}, y')$ ,  $ans_k$  is a valid solution for  $(I_k)_{y_{k+1}}$ .
- Suppose  $f_k(g_{k,y_{k+1}}(ans_k)) = y_{k+1}$ , then  $f_{\geq k}(y_k) = f_{\geq (k+1)}(f_k(g_{k,y_{k+1}}(ans_k))) = f_{\geq (k+1)}(y_{k+1})$ . However, the RHS is equal to  $y_d$  while the LHS is not equal to  $y_d$ . Therefore it must be the case that  $f_k(g_{k,y_{k+1}}(ans_k)) \neq y_{k+1}$ .

It follows that given a valid solution for  $(f', \{I'_y\}, \{g'_y\})$ , one can always find some  $k$  and a valid solution for  $(f_k, \{(I_k)_y\}, \{g_{k,y}\})$ . That is, there is a depth- $\gamma(O(\log \frac{N}{M}))$  reduction from solving  $\text{rPHP}(\mathcal{P})_{M \rightarrow N}$  to solving one of  $\{\text{rPHP}(\mathcal{P})_{(2^k M) \rightarrow (2^{k+1} M)}\}_{k \in [d]}$ . Composing this with the aforementioned depth-1 reduction from  $\text{rPHP}(\mathcal{P})_{M \rightarrow 2M}$  to  $\text{rPHP}(\mathcal{P})_{(2^k M) \rightarrow (2^{k+1} M)}$  completes our reduction.  $\square$

## B Comparing REFUTER with WRONGPROOF(Res)

We discuss the similarities and differences between the refuter problems and the WRONGPROOF problem. We first recall the formal definition of WRONGPROOF(Res) [BB17, GP18a]:

### Problem WRONGPROOF(Res)

Input: A CNF  $F$  with  $n$  variables and  $k$  clauses; a purported resolution refutation  $\Pi$  for  $F$  represented as  $C_{-k}, \dots, C_{-1}, C_0, C_1, \dots, C_{L-1}$ , where  $C_{-k}, \dots, C_{-1}$  are axioms of  $F$ ,  $C_{L-1} = \perp$ , and  $L = 2^{n^{\Omega(1)}}$ ; and a purported satisfying assignment  $\alpha \in \{0, 1\}^n$ .

Output: A number  $i \in [L]$  such that  $C_i$  is obtained by an invalid resolution derivation, or a number

$-k \leq j \leq -1$  such that  $\alpha$  does not satisfy  $C_j$ .

At first glance, the REFUTER problem looks similar to the WRONGPROOF problem. First, both problems take as input a purported (but not correct) resolution proof. Second, both are looking for an invalid derivation as a solution. Moreover, when we consider the resolution proof system (and consider refuting *width* lower bounds), both WRONGPROOF and REFUTER are PLS-complete.

However, we think that they are fundamentally different. One primary difference is the reason of totality: When introduced to a (non-promise) TFNP problem, the initial inquiry ought to be: *why is the problem total?* The totality of WRONGPROOF(Res) follows from the *reflection principle* for resolution [Pud20, BFI23], i.e., it is impossible to derive  $\perp$  from a satisfiable CNF. The same reasoning holds for every sound proof system, regardless of their power. However, the totality of REFUTER is far from trivial: *They rely on non-trivially proven width or size lower bounds.*

Furthermore, for comparison with REFUTER, we include a proof that WRONGPROOF(Res) is PLS-complete (this is a folklore result, see e.g., [BFI23]). The proof is seemingly similar to that of Theorem 5.1 and Theorem 6.1, but there are crucial differences. For example, the reduction from WRONGPROOF(Res) to PLS is uniform, since the totality of WRONGPROOF(Res) relies on simpler reasoning. In contrast, the uniform PLS-membership of REFUTER( $w(F \vdash_{\text{Res}} \perp)$ ) crucially relies on nice properties of the family of CNFs (e.g., EPHP), and it is possible that for some families, the refuter problem cannot be uniformly reduced to PLS at all. This demonstrates another difference between WRONGPROOF and REFUTER.

**Lemma B.1.** WRONGPROOF(Res) is in PLS.

*Proof.* Let  $(C_{-k}, \dots, C_{-1}, C_0, \dots, C_{L-1})$  be a purported resolution refutation of a CNF  $F$ , and  $\alpha$  be a purported satisfying assignment of  $F$ . We will reduce this WRONGPROOF(Res) instance to an instance  $S : \{-k, \dots, L-1\} \rightarrow \{-k, \dots, L-1\}$  of reversed ITER.

It would be convenient to think of a clause  $C_i$  as “active” if  $C_i(\alpha) = 0$ . An invalid derivation in the resolution refutation corresponds to an edge from an active node to an inactive node. For every  $i \in \{-k, \dots, L-1\}$ , if  $C_i(\alpha) = 1$  (i.e.,  $C_i$  is inactive), then we define  $S(i) = i$ . Otherwise, if  $i < 0$  (i.e.,  $C_i$  is an axiom *not* satisfied by  $\alpha$ ), then we define  $S(i) = 0$ , making  $i$  a solution since  $S(i) > i$ . Otherwise, suppose  $C_i$  is derived from  $C_j$  (i.e.,  $C_i$  is a *weakening* of  $C_j$ , or  $C_i$  is *resolved* from  $C_j$  and some other  $C_k$ ), then we define  $S(i) = j$ . If  $i$  is a solution for the reversed ITER instance  $S$ , then either  $j < i$  or  $j$  is inactive (which means  $S(j) = j$ ), and in either case  $i$  is a valid solution for WRONGPROOF(Res).  $\square$

*Remark 7.* The proof above is easy, but one can see that the crucial components are 1) the resolution proof system is sound, and 2) a resolution proof is a “DAG”-like structure. This proof strategy can potentially be easily extended to other proof systems with similar properties.

**Lemma B.2.** WRONGPROOF(Res) is PLS-hard.

*Proof.* We will reduce any reversed ITER instance to an instance of WRONGPROOF(Res). The construction below is very similar to the proof of Theorem 5.1. In fact, all clauses and derivations in the construction of the proof of Theorem 5.1 are sound except for the solutions of the given reversed ITER instance.

Let  $F$  be any satisfiable CNF with  $k$  clauses  $C_{-k}, \dots, C_{-1}$  and  $\alpha$  be any satisfying assignment of  $F$ . Without loss of generality assume there are two clauses  $C_{-2}$  and  $C_{-1}$  that we can apply a valid resolution step and call the resolved clause  $D$ . Let  $S : [L] \rightarrow [L]$  be an instance of reversed ITER where  $S(L) < L$ . We construct a purported resolution refutation  $\Pi = (C_{-k}, \dots, C_{-1}, C_0, \dots, C_{L-1})$  as follows:

- For every  $i$  such that  $S(i) = i$ , we let  $C_i := D$  to be *resolved* from  $C_{-2}$  and  $C_{-1}$ .
- For every  $i$  that is a solution for  $S$ , let  $C_i := \perp$  be a *weakening* from an axiom (say  $C_{-k}$ ). Note that this weakening step is invalid and  $C_i$  becomes a solution for the WRONGPROOF(Res) instance.

- Finally, for every  $i$  such that  $S(i) < i$  and  $S(S(i)) < S(i)$ , let  $C_i := \perp$  be a *weakening* of  $C_{S(i)}$ . Note that  $C_{S(i)}$  is also  $\perp$ , hence this is a valid derivation.

It is easy to see that the invalid derivations in  $\Pi$  correspond exactly to the solutions of  $S$ . □

## C Prover-Delayer Games, PLS, and the Proof of Lemma 7.10

In Section C.1, we provide a self-contained description of the transformation from a PLS formulation to a low-width resolution proof using *Prover-Delayer* game, along with several properties of this transformation that are useful when proving Lemma 7.10. We then prove Lemma 7.10 in Section C.2.

### C.1 From PLS to Resolution using Prover-Delayer Game

Introduced by Pudlák [Pud00], the *Prover-Delayer* game provides an elegant characterization of resolution width. There are two players in the game, the *Prover* (she) and the *Delayer* (he). Fixing an unsatisfiable CNF formula  $F$ , and let  $x = (x_1, \dots, x_n)$  be the variables in  $F$ . At first, the Prover’s memory is empty. Then, in each step, she can either

- *query* the Delayer for the value of a certain variable, and add that value to her memory;
- *forget* the value of a certain variable stored in her memory; or
- *output* a clause of  $F$  that is falsified by the partial assignment stored in her memory, which means she wins the game.

We assume the Delayer also has access to Prover’s memory. If the Prover queries a variable that is currently in its memory, then the Delayer’s answer must be consistent with the memory; otherwise, his answer could be arbitrary. Note that if the Prover queries a variable, forgets it, and queries it again, the Delayer is allowed to answer different values to these two queries of the same variable.

Of course, Prover can always win the game by querying all variables without forgetting any of them. However, for the connection with resolution width, her goal is to win the game with the minimum memory size, where the memory size is the maximum number of variables she remembered during the whole execution of the game. The Delayer is *adversarial* to Prover’s goal, i.e., wants her to spend as much memory as possible.

The following theorem shows that the minimum resolution width of an unsatisfiable CNF is characterized by the minimum Prover memory in the corresponding Prover-Delayer game.

**Theorem C.1** ([Pud00]). *For any unsatisfiable CNF formula  $F$ , there exists a width- $w$  resolution refutation of  $F$  if and only if there is a winning strategy for Prover using memory size  $w$  in the Prover-Delayer game for  $F$ .*

In this section, we prove the “if” direction in the previous lemma and highlight some nice properties of the obtained low-width resolution proof that will be helpful when proving Lemma 7.10.

**Making Prover’s strategy uniform.** Note that both Prover’s and Delayer’s strategies could be quite non-uniform in the Prover-Delayer game model described above. Here, we twist the model a little bit by allowing the Prover to explicitly store several *state registers* in its memory, besides a partial assignment of variables. These internal state registers are also counted in the memory size of her strategy. Later, in the proof of Lemma 7.10, it is more convenient to describe a uniform Prover’s strategy with state registers.

**From PLS to Prover-Delayer game.** A PLS formulation of a search problem  $\text{Search}(F_n)$  is a decision tree reduction  $(f_i, g_o)_{i,o \in M}$  from  $\text{Search}(F_n)$  to  $\text{ITER}_M$ . Let  $x = (x_1, \dots, x_n)$  be the variables in  $F_n$ , then we have  $S(v) := f_v(x)$ , where  $S$  is the successor function in the  $\text{ITER}_M$  instance reduced from  $\text{Search}(F_n)$ .

**Lemma C.2** (Folklore). *Given a PLS formulation  $(f, g)_{i,o \in M}$  of depth  $d$  for  $\text{Search}(F_n)$ , there exists a Prover's strategy of memory size  $O(d + \log M)$  for  $F_n$ .*

*Proof.* We say the Prover queries a decision tree  $T$  if she evaluates  $T(x)$ , and stores the queried variables in her memory in each step. Now we describe the Prover's strategy in what follows.

1. The Prover starts from the node 0 of the  $\text{ITER}_M$  instance and queries the decision tree  $f_0$ . If  $f_0(x) = 0$ , then 0 is a valid solution for the  $\text{ITER}_M$  instance, hence she can query  $g_0(x)$  to obtain a falsified clause in  $F_n$ . Otherwise, we say that the Prover is currently at node  $v = f_0$  and previously visited node 0.
2. Assume the Prover is at node  $v \in [M]$ , and the previous node she visited is  $u$ . She queries the decision tree  $f_v$  and obtains the next node  $w = f_v(x)$ .
- 3a. If  $w \leq v$ , then she has found a solution of the  $\text{ITER}_M$  instance. In particular, if  $w < v$  then the solution is  $v$ ; if  $w = v$  then the solution is  $u$ . Note that all the variables queried by  $f_u, f_v, g_u$  are still in her memory. Therefore, the clause  $F$  returned by  $g_v(x)$  (if  $w < v$ ) or  $g_u(x)$  (if  $w = v$ ) must be falsified by the variables in her memory.
- 3b. If  $w > v$ , then the Prover forgets all the variables that are queried in  $f_u$  and not queried in  $f_v$ . She then updates the current node as  $w$  and the previous node as  $v$ , and loops back to **Step 2**.

The Prover's strategy will always end, since the index of  $v$  increases in every step. The Prover needs to remember at most  $3d$  variables at any time, and  $O(d + \log M)$  bits to remember the current state to execute this strategy.  $\square$

By further examining the proof of [Lemma C.2](#), we obtain several properties that are useful for the proof of [Lemma 7.10](#) later.

**Observation C.3.** *In [Lemma C.2](#), the Prover's strategy can be implemented in a uniform manner if the PLS formulation  $(f, g)$  is given via oracle access.*

**Lemma C.4.** *In [Lemma C.2](#), there exists an efficient binary encoding of Prover's memory, such that:*

1. *The encoding has bit-length  $\text{poly}(d, \log M)$ .*
2. *It is computationally efficient to transform Prover's memory into an encoding and vice versa.*
3. *The encoding of the Prover's memory is lexicographically increasing as the Prover's strategy proceeds.*
4. *We say an encoding is invalid if it is in the wrong format, or Prover's internal state registers are inconsistent with the partial assignment of variables w.r.t.  $(f, g)$ . There is an efficient uniform algorithm for checking whether an encoding is invalid given oracle access to  $(f, g)$ .*

*Proof.* The Prover's internal state registers should store the index of the current node, the previous node, and its current location in the decision tree it is querying. This part is lexicographically increasing since we always have  $w > v$  in **Step 3b**. Our encoding consists of these internal state registers followed by the partial assignment of variables.

It is trivial to construct such an encoding scheme satisfying conditions 1,2,3, and easy to check the correctness of its format. To check the validity of an encoding, we can query the (at most 3) decision tree paths corresponding to the Prover's internal state registers, and check whether they are consistent with the partial assignment.  $\square$

## From Prover-Delayer games to resolution proofs

**Lemma C.5** ([Pud00]). *Given a Prover's strategy of memory cost  $w$  for  $F_n$ , there exists a width- $w$  resolution proof refuting  $F_n$ .*

*Proof.* Without loss of generality, we assume the Prover will not query a variable that is already in her memory.

We simulate the Prover's strategy. Initially, there is no variable stored in her memory, and it corresponds to the empty clause  $\perp$  at the end of the resolution proof. We then generate the resolution proof recursively by maintaining the Prover's memory and the current node of the resolution proof.

In each step, let  $\rho$  be the current partial assignment of variables stored in the Prover's memory. Define  $C(\rho)$  to be the only clause that is falsified by  $\rho$ , using only the variables that are set in  $\rho$ . For example, if  $\rho$  is  $\{x_1 = 1, x_2 = 0\}$ , the  $C(\rho) = \neg x_1 \vee x_2$ . The procedure will guarantee that  $C(\rho)$  is the clause of the current node in the resolution proof.

The Prover has three possible actions given  $\rho$  and its internal state register:

**FORGET** If the Prover decides to forget  $x_i$ , then we generate a new node  $C'$  with clause  $C(\rho_{-i})$ , where  $\rho_{-i}$  is the partial assignment by *forgetting* the value of  $x_i$  from  $\rho$ . We mark that the current node is derived by a *weakening* step from node  $C'$ . We then update the Prover's memory and continue our process at node  $C'$  recursively.

**QUERY** If the Prover queries  $x_i$ , then we generate two new nodes  $C^0, C^1$  with clauses  $C(\rho) \vee x_i$  and  $C(\rho) \vee \neg x_i$  respectively. We mark that the current node is derived by a *resolution* step from node  $C^0$  and  $C^1$ . We first recursively proceed to  $C^0$  by updating Prover's memory with  $x_i = 0$ , and then proceed to  $C^1$  with  $x_i = 1$ .

**OUTPUT** If the Prover outputs a falsified clause  $D$ , it must be the case that  $D$  is a sub-clause of  $C(\rho)$ . If  $D$  is equal to  $C(\rho)$ , then we simply stop; otherwise, we add a new node for the clause  $D$  and add one or more intermediate *weakening* steps towards the current node.

It is easy to verify that during the process,  $C(\rho)$  is always the clause of the current node in the resolution proof. This process will stop since the Prover will stop, and its correctness is guaranteed by the Prover's correctness. Finally, note that the largest clause ever generated in the resolution proof is upper bounded by the memory size of the Prover.  $\square$

### C.2 Proof of Lemma 7.10

**Lemma 7.10.** *For any family of unsatisfiable CNF  $\mathcal{F}$  that has no  $\text{polylog}(N)$ -width resolution refutation,*

$$\text{REFUTER}_{d,M}(\text{Search}(\mathcal{F}) \rightarrow \text{ITER}) \leq_m \text{REFUTER}(w(\mathcal{F} \vdash_{\text{Res}} \perp) < w_0)$$

for some  $w_0 = \text{polylog}(N)$  that may depend on  $d, M$ .

Furthermore, this reduction is uniform when  $\mathcal{F}$  is a uniform family of unsatisfiable CNFs.

*Proof.* Let  $\mathcal{F} = \{F_n\}_{n \in \mathbb{N}}$ . Suppose we are given a purported depth- $d$  reduction  $(f, g)$  from  $\text{Search}(F_n)$  to  $\text{ITER}_M$ , i.e.,  $(f, g)$  is a PLS formulation for  $\text{Search}(F_n)$ . We now use  $(f, g)$  to construct a purported resolution refutation  $(C_0, \dots, C_{L-1})$  for  $F_n$  with width  $w_0 = \text{poly}(d, \log M)$ , while satisfying the following two conditions.

1. Given any index  $i \in [L]$ , the  $i$ -th node  $C_i$  can be calculated in  $\text{polylog}(n)$  queries to  $(f, g)$  uniformly.
2. If the  $i$ -th node is invalid, one can recover a pair  $(\rho, \sigma^*)$  that refutes  $(f, g)$  uniformly given  $i$ .

As a high-level plan, we first apply the procedures described in [Lemma C.2](#) which transforms the PLS formulation  $(f, g)$  to a Prover's strategy of  $O(d + \log M)$  memory for the Prover-Delayer game. We then use the procedure described in [Lemma C.5](#) to convert such a strategy into a width- $O(d + \log M)$  resolution proof for  $F_n$ .

We now specify the details in these two steps to make sure the two conditions are met. The first condition can be achieved by letting an index  $i \in [L]$  to be an encoding of the Prover's memory in a single step, as described in [Lemma C.4](#).<sup>34</sup>

By [Observation C.3](#), given any valid index (encoding), one can calculate the next action of the Prover using  $O(d)$  number of queries to  $(f, g)$ . Then, for all three actions **{FORGET, QUERY, OUTPUT}**, we can also calculate the indices of the one or two previous nodes in the resolution proof uniformly. For any invalid index (encoding), we pad a trivially correct resolution node using axioms from the beginning.

To see the second condition above, note that the procedure described in [Lemma C.5](#) generates an invalid node of the resolution proof only when the Prover is taking an **OUTPUT** step. That is, when the Prover outputs a falsified clause  $D$ ,  $D$  might not be a sub-clause of  $C(\rho)$ , so the *weakening* steps from  $D$  to  $C(\rho)$  will be wrong. This happens because a solution  $o$  of the  $\text{ITER}_M$  instance is found by the Prover, but querying  $g_o$  does not lead to a clause  $D$  in  $F_n$  that is falsified by the current partial assignment  $\rho$  in her memory.

Note that the partial assignment  $\rho$  and  $o^*$  used to refute  $(f, g)$  can be recovered from the index (encoding) of the invalid node. We then complement the partial assignment  $\rho$  to  $\rho'$  by setting all unassigned variables that appear in  $D$  to satisfy clause  $D$ . By definition, the pair  $(\rho', o)$  is a valid solution to refute the reduction  $(f, g)$ .

Finally, if  $\mathcal{F}$  is a uniform family of formulas, then the whole reduction is also uniform. □

---

<sup>34</sup>Note that the encoding described in [Lemma C.4](#) is lexicographically increasing, but we can easily make it lexicographically decreasing by flipping all the bits.