

Learning Algorithms Versus Automatability of Frege Systems

Ján Pich ✉

University of Oxford, UK

Rahul Santhanam ✉

University of Oxford, UK

Abstract

We connect learning algorithms and algorithms automating proof search in propositional proof systems: for every sufficiently strong, well-behaved propositional proof system P , we prove that the following statements are equivalent,

- **Provable learning.** P proves efficiently that p -size circuits are learnable by subexponential-size circuits over the uniform distribution with membership queries.
- **Provable automatability.** P proves efficiently that P is automatable by non-uniform circuits on propositional formulas expressing p -size circuit lower bounds.

Here, P is sufficiently strong and well-behaved if I.-III. holds: I. P p -simulates Jeřábek's system WF (which strengthens the Extended Frege system EF by a surjective weak pigeonhole principle); II. P satisfies some basic properties of standard proof systems which p -simulate WF; III. P proves efficiently for some Boolean function h that h is hard on average for circuits of subexponential size. For example, if III. holds for $P = \text{WF}$, then Items 1 and 2 are equivalent for $P = \text{WF}$. The notion of automatability in Item 2 is slightly modified so that the automating algorithm outputs a proof of a given formula (expressing a p -size circuit lower bound) in p -time in the length of the shortest proof of a closely related but different formula (expressing an average-case subexponential-size circuit lower bound).

If there is a function $h \in \text{NE} \cap \text{coNE}$ which is hard on average for circuits of size $2^{n/4}$, for each sufficiently big n , then there is an explicit propositional proof system P satisfying properties I.-III., i.e. the equivalence of Items 1 and 2 holds for P .

2012 ACM Subject Classification Theory of computation

Keywords and phrases learning algorithms, automatability, proof complexity

Digital Object Identifier 10.4230/LIPIcs.ICALP.2022.101

Category Track A: Algorithms, Complexity and Games

Related Version *Previous Version:* <https://arxiv.org/abs/2111.10626>

Funding *Ján Pich:* This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 890220. *Ján Pich* received support from the Royal Society University Research Fellowship URF\R1\211106. *Rahul Santhanam:* Rahul Santhanam is partially funded by the EPSRC New Horizons grant EP/V048201/1: "Structure versus Randomness in Algorithms and Computation".



Acknowledgements We would like to thank Moritz Müller, Jan Krajíček, Iddo Zameret and anonymous reviewers for comments on a draft of the paper.

1 Introduction

Learning algorithms and automatability algorithms searching for proofs in propositional proof systems are central concepts in complexity theory, but a priori they appear rather unrelated.



© Ján Pich and Rahul Santhanam;

licensed under Creative Commons License CC-BY 4.0

49th International Colloquium on Automata, Languages, and Programming (ICALP 2022).

Editors: Mikołaj Bojańczyk, Emanuela Merelli, and David P. Woodruff;

Article No. 101; pp. 101:1–101:20



Leibniz International Proceedings in Informatics

LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



Learning algorithms. In the PAC model of learning introduced by Valiant [34], a circuit class \mathcal{C} is learnable by a randomized algorithm L over the uniform distribution, up to error ϵ , with confidence δ and membership queries, if for every Boolean function f computable by a circuit from \mathcal{C} , when given oracle access to f , L outputs with probability $\geq \delta$ over the uniform distribution a circuit computing f on $\geq (1 - \epsilon)$ inputs. An important task of learning theory is to find out if standard circuit classes such as P/poly are learnable by efficient circuits. A way to approach the question is to connect the existence of efficient learning algorithms to other standard conjectures in complexity theory. For example, we can try to prove that efficient learning of P/poly is equivalent to $P = NP$ or to the non-existence of strong pseudorandom generators. In both cases one implication is known: $P = NP$ implies efficient learning of P/poly (with small error and high confidence) which in turn breaks pseudorandom generators. However, while some progress on the opposite implications has been made, they remain open, cf. [2, 33].

Automatability. The notion of automatability was introduced in the work of Bonet, Pitassi and Raz [6]. A propositional proof system P is automatable if there is an algorithm A such that for every tautology ϕ , A finds a P -proof of ϕ in p-time in the size of the shortest P -proof of ϕ . That is, even if P does not prove all tautologies efficiently, it can still be automatable. Establishing (non-)automatability results for concrete proof systems is one of the main tasks of proof complexity. This led to many attempts to link the notion of automatability to other standard complexity-theoretic conjectures. For example, recently Atserias and Müller [3] proved that automating Resolution is NP-hard and their work has been extended to other weak proof systems, e.g. [12, 13, 14]. For stronger systems, it is known that automating Extended Frege system EF, Frege or even constant-depth Frege would break specific cryptographic assumptions such as the security of RSA or Diffie-Hellman scheme, cf. [23, 6, 5]. It remains, however, open to obtain non-automatability of strong systems like Frege under a generic assumption such as the existence of strong pseudorandom generators, let alone to prove the equivalence between such notions.

In the present paper we derive a conditional equivalence between learning algorithms for p-size circuits and automatability of proof systems on tautologies encoding circuit lower bounds.

1.1 Our result

An ideal connection between learning and automatability would say that for standard proof systems P ,

“ P is automatable if and only if P/poly is learnable efficiently”.

We establish this modulo some provability conditions and a change of parameters. Additionally, we need to consider automatability only w.r.t. formulas encoding circuit lower bounds. More precisely, denote by $\text{tt}(f, s)$ a propositional formula which expresses that boolean function f represented by its truth-table is not computable by a boolean circuit of size s represented by free variables, see Section 3. So $\text{tt}(f, s)$ is a tautology if and only if f is hard for circuits of size s . Note that f is represented by 2^n bits, if n is the number of inputs of f , so the size of $\text{tt}(f, s)$ is $2^{O(n)}$. Similarly, let $\text{tt}(f, s, t)$ be a formula expressing that circuits of size s fail to compute f on $\geq t$ -fraction of inputs. In our main result (Theorem 1)

we use a slightly modified notion of automatability where the automating algorithm for a proof system P is non-uniform and outputs a P -proof of a given formula $\text{tt}(f, n^{O(1)})$ in p -time in the size of the shortest P -proof of $\text{tt}(f, 2^{n^{o(1)}})$, $1/2 - 1/2^{n^{o(1)}}$, see Section 3.¹

► **Theorem 1** (Informal, cf. Theorem 18). *Let P be any propositional proof system which APC_1 -provably p -simulates WF and satisfies some basic properties, e.g. $P = \text{WF}$. Moreover, assume that P proves efficiently $\text{tt}(h, 2^{n/4}, 1/2 - 1/2^{n/4})$ for some boolean function h . Then, the following statements are equivalent:*

1. **Provable learning.** *P proves efficiently that p -size circuits are learnable by $2^{n^{o(1)}}$ -size circuits, over the uniform distribution, up to error $1/2 - 1/2^{n^{o(1)}}$, with membership queries and confidence $1/2^{n^{o(1)}}$.*
2. **Provable automatability.** *P proves efficiently that P is automatable by non-uniform circuits on formulas $\text{tt}(f, n^{O(1)})$.*

WF is an elegant strengthening of EF introduced by Jeřábek [15], which corresponds to the theory of approximate counting APC_1 , a theory formalizing probabilistic p -time reasoning, see Section 2.2. Concrete proof systems which APC_1 -provably p -simulate WF and satisfy the basic properties from Theorem 1 include WF itself or even much stronger systems such as set theory ZFC (if we interpret ZFC as a suitable system for proving tautologies, see Section 5). We emphasize that the conditional equivalence from Theorem 1 holds for any sufficiently strong proof system satisfying some basic properties. The error and confidence of learning algorithms can be amplified “for free”, see Section 2.1, but we did not make the attempts to prove that the amplification is efficiently provable already in WF .

Perhaps the most unusual aspect of Theorem 1 is its usage of metamathematics: we do not prove the equivalence between automatability and learning but between *provable* automatability and *provable* learning. We believe that the usage of metamathematics is not a substantial deviation. It would be very surprising if, e.g., an efficient learning algorithm existed but not provably (in some natural proof system).

Plausibility of the assumption. The main assumption in Theorem 1 is the provability of a circuit lower bound $\text{tt}(h, 2^{n/4}, 1/2 - 1/2^{n/4})$. This assumption has an interesting status. Razborov’s conjecture about hardness of Nisan-Wigderson generators implies a conditional hardness of formulas $\text{tt}(h, n^{O(1)})$ for Frege (for every h), cf. [31], and it is possible to consider extensions of the conjecture to all standard proof systems, even set theory ZFC . On the other hand, all major circuit lower bounds for weak circuit classes and explicit boolean functions are known to be efficiently provable in EF ², cf. [29, 25]. If we believe that explicit circuit lower bounds such as $\text{tt}(h, 2^{n/4}, 1/2 - 1/2^{n/4})$, for some $h \in \text{EXP}$, are true, it is also perfectly plausible that they are efficiently provable in a standard proof system such as ZFC ³ or EF . Notably, if EF proves efficiently $\text{tt}(h, 2^{n/4})$ for some boolean function h , then EF simulates WF , cf. [22, Lemma 19.5.4]. If there is a p -time algorithm which given a string of length 2^n (specifying the size of $\text{tt}(h, 2^{n/4})$) generates an EF -proof of $\text{tt}(h, 2^{n/4})$, then EF is p -equivalent to WF . To see that, combine Lemma 12 with the fact (proved in [15]) that APC_1 proves the reflection principle for WF .

¹ We believe that the gap between $\text{tt}(f, n^{O(1)})$ and $\text{tt}(f, 2^{n^{o(1)}})$, $1/2 - 1/2^{n^{o(1)}}$ can be almost closed, if one uses learning of subexponential-size circuits instead of p -size circuits in Item 1 of Theorem 1 and tt -formulas expressing subexponential-size circuit lower bounds in Item 2.

² This has not been verified for lower bounds obtained via the algorithmic method of Williams [35].

³ Efficient provability of $\text{tt}(h, 2^{n/4}, 1/2 - 1/2^{n/4})$ in ZFC , for some $h \in \text{EXP}$, would follow from the standard provability of this lower bound in ZFC .

As a corollary of Theorem 1 we show that, under a standard hardness assumption, there is an explicit proof system P for which the equivalence holds. This, follows, essentially, by “hard-wiring” tautologies $\text{tt}(h, 2^{n/4}, 1/2 - 1/2^{n/4})$ to WF.

► **Corollary 2** (cf. Corollary 22). *Assume there is a $\text{NE} \cap \text{coNE}$ -function $h_n : \{0, 1\}^n \mapsto \{0, 1\}$ such that for each sufficiently big n , h_n is not $(1/2 + 1/2^{n/4})$ -approximable by $2^{n/4}$ -size circuits.⁴ Then there is a proof system P (which can be described explicitly given the definition of h_n) such that Items 1 and 2 from Theorem 1 are equivalent.*

The proof of Theorem 1 reveals also a proof complexity collapse which we discuss in the arXiv version of the paper.

1.2 Outline of the proof

Our starting point for the derivation of Theorem 1 is a relation between natural proofs and automatability which goes back to a work of Razborov and Krajíček. Razborov [30, 28] proved that certain theories of bounded arithmetic cannot prove explicit circuit lower bounds assuming strong pseudorandom generators exist. Krajíček [19, 21] developed the concept of feasible interpolation (a weaker version of automatability, cf. [22]) and reformulated Razborov’s unprovability result in this language, see [22, Section 17.9] for more historical remarks.

► **Theorem 3** (Razborov-Krajíček [30, 28, 19] - informal version). *Let P be a proof system which simulates EF and satisfies some basic properties. If P is automatable and P proves efficiently $\text{tt}(h, n^{O(1)})$ for some function h , then there are P/poly-natural proofs useful against P/poly.*

The second crucial ingredient we will use is a result of Carmosino, Impagliazzo, Kabanets and Kolokolova, who showed that natural proofs can be turned into learning algorithms [8]. This allows us to conclude the following.

► **Theorem 4** (Informal). *Let P be a proof system simulating EF and satisfying some basic properties. If P proves efficiently $\text{tt}(h, n^{O(1)})$ for some function h , then automatability of P implies the existence of subexponential-size circuits learning p -size circuits over the uniform distribution, with membership queries.*

Theorem 4 directly implies that if strong pseudorandom generators exist and EF proves efficiently $\text{tt}(h, n^{O(1)})$ for some h , then EF is automatable if and only if there are subexponential-size circuits learning p -size circuits over the uniform distribution, with membership queries. The disadvantage of this observation is that, unlike in Theorem 1, its assumptions are known to imply that both sides of the desired equivalence are false.

We note that the proof of Theorem 4 can be used to show also that optimal and automatable proof systems imply learning algorithms. In fact, it is possible to prove, unconditionally, that there is some propositional proof system P such that automatability of P is equivalent to the existence of subexponential-size circuits infinitely often learning P/poly over the uniform distribution. The proof is, however, non-constructive so (unlike in Corollary 2) we do not know which system P satisfies the equivalence. We discuss these results in more detail in the arXiv version of the paper.

⁴ A circuit C with n inputs γ -approximates function $f : \{0, 1\}^n \mapsto \{0, 1\}$ if $\Pr_{x \in \{0, 1\}^n} [C(x) = f(x)] \geq \gamma$.

The entrance of metamathematics. Unfortunately, it is unclear how to derive the opposite implication in Theorem 4. We do not know how to automate, say, EF assuming just the existence of efficient learning algorithms. In order to get the reverse, we need to assume that an efficient learning algorithm is provably correct in a proof system P , which p -simulates WF. For simplicity, let $P = \text{WF}$. If we assume that WF proves efficiently for some small circuits that they can learn p -size circuits, we can show that there are small circuits such that WF proves efficiently that these circuits automate WF on formulas $\text{tt}(f, n^{O(1)})$. In more detail, we first formalize in APC_1 the implication that WF-provable learning yields automatability of WF on $\text{tt}(f, n^{O(1)})$ - if a learning circuit A does not find a small circuit for a given function f , the automating circuit uses WF-proof of the correctness of A to produce a short WF-proof of $\text{tt}(f, n^{O(1)})$. Then, we translate the APC_1 -proof to WF and conclude that WF proves that WF-provable learning implies automatability of WF. This allows us to show that if we have WF-provable learning, then WF is WF-provably automatable on $\text{tt}(f, n^{O(1)})$.

It is important that assuming WF-provable learning, we are able to derive WF-provable automatability of WF, and not just automatability of WF. This makes it possible to obtain the opposite direction and establish the desired equivalence: If we know that WF proves that WF is automatable, we can formalize the proof of Theorem 4 in WF and conclude the existence of WF-provable learning algorithms.

Benefits of bounded arithmetic. The proof of Theorem 1 relies heavily on formalizations. Among other things we need to formalize the result of Carmosino, Impagliazzo, Kabanets and Kolokolova in APC_1^5 , and use an elaborated way of expressing complex statements about metacomplexity by propositional formulas: existential quantifiers often need to be witnessed before translating them to propositional setting. The framework of bounded arithmetic allows us to deal with these complications in an elegant way: we often reason in bounded arithmetic, possibly using statements of higher quantifier complexity, and only subsequently translate the outcomes to propositional logic, if the resulting (proved) statement has coNP form. Notably, already propositional formulas expressing probabilities in the definition of learning algorithms require more advanced tools - the probabilities are encoded using suitable Nisan-Wigderson generators which come out of the notion of approximate counting in APC_1 , cf. Section 3.2.

1.3 Related results

Learning algorithms and automatability have been linked already in the work of Alekhovich, Braverman, Feldman, Klivans and Pitassi [1], who showed an informal connection between learning of weak circuit classes and automatability of some weak systems such as tree-like Resolution. As already mentioned, Atserias and Müller [3] proved that automating Resolution is NP-hard and their work has been extended to other weak proof systems, see e.g. [12, 13, 14]. A direct consequence of these results is that efficient algorithms automating the respective proof systems can be used to learn efficiently classes like P/poly . A major difference between these results and ours is that for our results to apply, the proof system needs to be sufficiently strong, while for the other results, the proof system needs to be weak (in the sense that lower bounds for the system are already known).

⁵ We will actually formalize “CIKK” just conditionally, in order to avoid the formalization of Bertrand’s postulate.

1.4 Open problems

Unconditional equivalence between learning and automatability. Is it possible to avoid the assumption on the provability of a circuit lower bound in Theorem 1 and establish an unconditional equivalence between learning and automatability?

Complexity theory from the perspective of metamathematics. Our results demonstrate that in the context of metamathematics it is possible to establish some complexity-theoretic connections which we are not able to establish otherwise. We exploit the metamathematical nature of the notion of automatability: efficient P -provability of the correctness of an algorithm implies efficient P -provability of automatability of P . Is it possible to take advantage of metamathematics in other contexts and resolve other important open problems in this setting? For example, could we get a version of the desired equivalence between the existence of efficient learning algorithms and the non-existence of cryptographic pseudorandom generators, cf. [26, 33, 27]? The question of basing cryptography on a worst-case assumption such as $P \neq NP$ could be addressed in this setting by showing that if a sufficiently strong proof system P proves efficiently that there is no strong pseudorandom generator⁶, then P is p-bounded.

Circuit lower bound tautologies. How essential are circuit lower bound tautologies in our results? Consider fundamental questions of proof complexity (p-boundness, optimality, automatability) w.r.t. formulas $\text{tt}(f, s)$. Do they coincide with the original ones? Are formulas $\text{tt}(f, s)$ the hardest ones, do they admit optimal proof systems, or can we turn automatability on formulas $\text{tt}(f, s)$ into automatability on all formulas?

2 Preliminaries

2.1 Natural proofs and learning algorithms

$[n]$ denotes $\{1, \dots, n\}$. $\text{Circuit}[s]$ denotes fan-in two Boolean circuits of size at most s . The size of a circuit is the number of gates.

► **Definition 5** (Natural property [32]). *Let $m = 2^n$ and $s, d : \mathbb{N} \mapsto \mathbb{N}$. A sequence of circuits $\{C_{2^n}\}_{n=1}^\infty$ is a $\text{Circuit}[s(m)]$ -natural property useful against $\text{Circuit}[d(n)]$ if*

1. Constructivity. C_m has m inputs and size $s(m)$,
2. Largeness. $\Pr_x[C_m(x) = 1] \geq 1/m^{O(1)}$,
3. Usefulness. *For each sufficiently big m , $C_m(x) = 1$ implies that x is a truth-table of a function on n inputs which is not computable by circuits of size $d(n)$.*

► **Definition 6** (PAC learning). *A circuit class \mathcal{C} is learnable over the uniform distribution by a circuit class \mathcal{D} up to error ϵ with confidence δ , if there are randomized oracle circuits L^f from \mathcal{D} such that for every Boolean function $f : \{0, 1\}^n \mapsto \{0, 1\}$ computable by a circuit from \mathcal{C} , when given oracle access to f , input 1^n and the internal randomness $w \in \{0, 1\}^*$, L^f outputs the description of a circuit satisfying*

$$\Pr_w[L^f(1^n, w) \text{ (1} - \epsilon)\text{-approximates } f] \geq \delta.$$

⁶ The formalization of this statement would assume the existence of a p-size circuit which for any p-size circuit defining a potential pseudorandom generator outputs its distinguisher.

L^f uses non-adaptive membership queries if the set of queries which L^f makes to the oracle does not depend on the answers to previous queries. L^f uses random examples if the set of queries which L^f makes to the oracle is chosen uniformly at random.

In this paper, PAC learning always refers to learning over the uniform distribution. While, a priori, learning over the uniform distribution might not reflect real-world scenarios very well (and on the opposite end, learning over all distributions is perhaps overly restrictive), as far as we can tell it is possible that PAC learning of p -size circuits over the uniform distribution implies PAC learning of p -size circuits over all p -samplable distributions. Binnendyk, Carmosino, Kolokolova, Ramyaa and Sabin [4] proved the implication, if the learning algorithm in the conclusion is allowed to depend on the p -samplable distribution.

Boosting confidence and reducing error. The confidence of the learner and its error can be improved generically, see the arXiv version of the paper. We can thus often ignore the optimisation of these parameters.

Natural proofs vs learning algorithms. Natural proofs are actually equivalent to efficient learning algorithms with suitable parameters. In this paper we need just one implication.

► **Theorem 7** (Carmosino-Impagliazzo-Kabanets-Kolokolova [8]). *Let R be a P/poly -natural property useful against $\text{Circuit}[n^k]$ for $k \geq 1$. Then, for each $\gamma \in (0, 1)$, $\text{Circuit}[n^{k\gamma/a}]$ is learnable by $\text{Circuit}[2^{O(n^\gamma)}]$ over the uniform distribution with non-adaptive membership queries, confidence 1, up to error $1/n^{k\gamma/a}$, where a is an absolute constant.*

2.2 Bounded arithmetic and propositional logic

Theories of bounded arithmetic capture various levels of feasible reasoning and present a uniform counterpart to propositional proof systems.

The first theory of bounded arithmetic formalizing p -time reasoning was introduced by Cook [10] as an equational theory PV . We work with its first-order conservative extension PV_1 from [24]. The language of PV_1 , denoted PV as well, consists of symbols for all p -time algorithms given by Cobham's characterization of p -time functions, cf. [9]. A PV -formula is a first-order formula in the language PV . Σ_0^b ($=\Pi_0^b$) denotes PV -formulas with only sharply bounded quantifiers $\exists x, x \leq |t|$, $\forall x, x \leq |t|$, where $|t|$ is “the length of the binary representation of t ”. Inductively, Σ_{i+1}^b resp. Π_{i+1}^b is the closure of Π_i^b resp. Σ_i^b under positive Boolean combinations, sharply bounded quantifiers, and bounded quantifiers $\exists x, x \leq t$ resp. $\forall x, x \leq t$. Predicates definable by Σ_i^b resp. Π_i^b formulas are in the Σ_i^p resp. Π_i^p level of the polynomial hierarchy, and vice versa. PV_1 is known to prove $\Sigma_0^b(PV)$ -induction:

$$A(0) \wedge \forall x (A(x) \rightarrow A(x+1)) \rightarrow \forall x A(x),$$

for Σ_0^b -formulas A , cf. Krajíček [18].

Buss [7] introduced the theory S_2^1 extending PV_1 with the Σ_1^b -length induction:

$$A(0) \wedge \forall x < |a|, (A(x) \rightarrow A(x+1)) \rightarrow A(|a|),$$

for $A \in \Sigma_1^b$. S_2^1 proves the sharply bounded collection scheme $BB(\Sigma_1^b)$:

$$\forall i < |a| \exists x < a, A(i, x) \rightarrow \exists w \forall i < |a|, A(i, [w]_i),$$

for $A \in \Sigma_1^b$ ($[w]_i$ is the i th element of the sequence coded by w), which is unprovable in PV_1 under a cryptographic assumption, cf. [11]. On the other hand, S_2^1 is $\forall\Sigma_1^b$ -conservative over PV_1 . This is a consequence of Buss's witnessing theorem stating that $S_2^1 \vdash \exists y, A(x, y)$ for $A \in \Sigma_1^b$ implies $PV_1 \vdash A(x, f(x))$ for some PV-function f .

Following a work by Krajíček [20], Jeřábek [15, 16, 17] systematically developed a theory APC_1 capturing probabilistic p-time reasoning by means of approximate counting.⁷ The theory APC_1 is defined as $PV_1 + dWPHP(PV)$ where $dWPHP(PV)$ stands for the dual (surjective) pigeonhole principle for PV-functions, i.e. for the set of all formulas

$$x > 0 \rightarrow \exists v < x(|y| + 1) \forall u < x|y|, f(u) \neq v,$$

where f is a PV-function which might involve other parameters not explicitly shown. We devote Section 2.3 to a more detailed description of the machinery of approximate counting in APC_1 .

Any Π_1^b -formula ϕ provable in PV_1 can be expressed as a sequence of tautologies $\|\phi\|_n$ with proofs in the Extended Frege system EF which are constructible in p-time (given a string of the length n), cf. [10]. Similarly, Π_1^b -formulas provable in APC_1 translate to tautologies with p-time constructible proofs in WF , an extension of EF introduced by Jeřábek [15]. We describe the translation and system WF in more detail below.

As it is often easier to present a proof in a theory of bounded arithmetic than in the corresponding propositional system, bounded arithmetic functions, so to speak, as a uniform language for propositional logic.

We refer to Krajíček [22] for basic notions in proof complexity.

► **Definition 8** (*WF (WPHP Frege)*, cf. Jeřábek [15]). *Let L be a finite and complete language for propositional logic, i.e. L consists of finitely many boolean connectives of constant arity such that each boolean function of every arity can be expressed by an L -formula, and let \mathcal{R} be a finite, sound and implicationally complete set of Frege rules (in the language L). A WF -proof of a (L -)circuit A is a sequence of circuits A_0, \dots, A_k such that $A_k = A$, and each A_i is derived from some $A_{j_1}, \dots, A_{j_\ell}$, $j_1, \dots, j_\ell < i$ by a Frege rule from \mathcal{R} , or it is similar to some A_j , $j < i$, or it is the $dWPHP$ axiom,*

$$\bigvee_{\ell=1}^m (r_\ell \neq C_{i,\ell}(D_{i,1}, \dots, D_{i,n})),$$

where $n < m$ and r_ℓ are pairwise distinct variables which do not occur in circuits A , $C_{i,\ell}$, or A_j for $j < i$, but may occur in circuits $D_{i,1}, \dots, D_{i,n}$.

The similarity rule in Definition 8 is verified by a specific p-time algorithm which checks that circuits A_i and A_j can be “unfolded” to the same (possible huge) formula, cf. [15, Lemma 2.2.]. Intuitively, the $NLOG (\subseteq P)$ algorithm recognizes if two circuits are not similar by guessing a partial path through them, going from the output to the inputs, where on at least one instruction the circuits disagree. As defined WF depends on the choice of Frege rules and language L , but for each choice the resulting systems are p-equivalent, so we can identify them. The $dWPHP$ axiom refers to “dual weak pigeonhole principle” postulating the existence of an element r_1, \dots, r_m outside the range of a p-size map $C_{i,1}, \dots, C_{i,m} : \{0, 1\}^n \mapsto \{0, 1\}^m$.

⁷ Krajíček [20] introduced a theory BT defined as $S_2^1 + dWPHP(PV)$ and proposed it as a theory for probabilistic p-time reasoning.

The dWPHP axiom comes with a specification of circuits $C_{i,1}, \dots, C_{i,m}, D_{i,1}, \dots, D_{i,n}$ so that we can recognize the axiom efficiently. The role of circuits $D_{i,1}, \dots, D_{i,n}$ in the dWPHP axiom is to allow WF to postulate not only that r_1, \dots, r_m is not the output of $C_{i,1}, \dots, C_{i,m}$ on a specific input x_1, \dots, x_n but to postulate that r_1, \dots, r_m is not the output of $C_{i,1}, \dots, C_{i,m}$ on other inputs (which could depend on r_1, \dots, r_m) either.

The translation of a Π_1^b formula ϕ into a sequence of propositional formulas $\|\phi\|_{\bar{n}}$ works as follows. For each PV-function $f(x_1, \dots, x_k)$ and numbers n_1, \dots, n_k we have a p-size circuit C_f computing the restriction $f : 2^{n_1} \times \dots \times 2^{n_k} \mapsto 2^{b(n_1, \dots, n_k)}$, where b is a suitable “bounding” polynomial for f . The formula $\|f\|_{\bar{n}}(p, q, r)$ expresses that C_f outputs r on input p , with q being the auxiliary variables corresponding to the nodes of C_f . The formula $\|\phi(x)\|_{\bar{n}}(p, q)$ is defined as $\|\phi'(x)\|_{\bar{n}}(p, q)$, where ϕ' is the negation normal form of ϕ , i.e. negations in ϕ' are only in front of atomic formulas. The formula $\|\phi'(x)\|_{\bar{n}}(p, q)$ is defined inductively in a straightforward way so that $\|\dots\|$ commutes with \vee, \wedge . The atoms p correspond to variables x , atoms q correspond to the universally quantified variables of ϕ and to the outputs and auxiliary variables of circuits C_f for functions f appearing in ϕ . Sharply bounded quantifiers are replaced by polynomially big conjunctions resp. disjunctions. For the atomic formulas we have,

$$\begin{aligned} \|f(x) = g(x)\|_{\bar{n}} &:= \|f(x)\|_{\bar{n}}(p, q, r) \wedge \|g(x)\|_{\bar{n}}(p, q', r') \rightarrow \bigwedge_i r_i = r'_i, \\ \|\neg f(x) = g(x)\|_{\bar{n}} &:= \|f(x)\|_{\bar{n}}(p, q, r) \wedge \|g(x)\|_{\bar{n}}(p, q', r') \rightarrow \neg \bigwedge_i r_i = r'_i, \\ \|f(x) \leq g(x)\|_{\bar{n}} &:= \|f(x)\|_{\bar{n}}(p, q, r) \wedge \|g(x)\|_{\bar{n}}(p, q', r') \rightarrow \bigwedge_i (r_i \wedge \bigwedge_{j>i} (r_j = r'_j) \rightarrow r'_i), \\ \|\neg f(x) \leq g(x)\|_{\bar{n}} &:= \|f(x)\|_{\bar{n}}(p, q, r) \wedge \|g(x)\|_{\bar{n}}(p, q', r') \rightarrow \neg \bigwedge_i (r_i \wedge \bigwedge_{j>i} (r_j = r'_j) \rightarrow r'_i). \end{aligned}$$

2.3 Approximate counting

In order to prove our results we will need to use Jeřábek’s theory of approximate counting. This section recalls the properties of APC_1 we will need.

By a definable set we mean a collection of numbers satisfying some formula, possibly with parameters. When a number a is used in a context which asks for a set it is assumed to represent the integer interval $[0, a)$, e.g. $X \subseteq a$ means that all elements of set X are less than a . If $X \subseteq a$, $Y \subseteq b$, then $X \times Y := \{bx + y \mid x \in X, y \in Y\} \subseteq ab$ and $X \dot{\cup} Y := X \cup \{y + a \mid y \in Y\} \subseteq a + b$. Rational numbers are assumed to be represented by pairs of integers in the natural way. We use the notation $x \in \text{Log} \leftrightarrow \exists y, x = |y|$ and $x \in \text{LogLog} \leftrightarrow \exists y, x = \|y\|$.

Let $C : 2^n \rightarrow 2^m$ be a circuit and $X \subseteq 2^n, Y \subseteq 2^m$ definable sets. We write $C : X \twoheadrightarrow Y$ if $\forall y \in Y \exists x \in X, C(x) = y$. Jeřábek [17] gives the following definitions in APC_1 (but they can be considered in weaker theories as well).

► **Definition 9.** Let $X, Y \subseteq 2^n$ be definable sets, and $\epsilon \leq 1$. The size of X is approximately less than the size of Y with error ϵ , written as $X \preceq_\epsilon Y$, if there exists a circuit C , and $v \neq 0$ such that

$$C : v \times (Y \dot{\cup} \epsilon 2^n) \twoheadrightarrow v \times X.$$

$X \approx_\epsilon Y$ stands for $X \preceq_\epsilon Y$ and $Y \preceq_\epsilon X$.

101:10 Learning Algorithms Versus Automatability of Frege Systems

Since a number s is identified with the interval $[0, s)$, $X \preceq_\epsilon s$ means that the size of X is at most s with error ϵ .

The definition of $X \preceq_\epsilon Y$ is an unbounded $\exists\Pi_2^b$ -formula even if X, Y are defined by circuits so it cannot be used freely in bounded induction. Jeřábek [17] solved this problem by working in HARD^A , a conservative extension of APC_1 , defined as a relativized theory $\text{PV}_1(\alpha) + dWPHP(\text{PV}(\alpha))$ extended with axioms postulating that $\alpha(x)$ is a truth-table of a function on $\|x\|$ variables hard on average for circuits of size $2^{\|x\|/4}$, see Section 3.2. In HARD^A there is a $\text{PV}_1(\alpha)$ function Size approximating the size of any set $X \subseteq 2^n$ defined by a circuit C so that $X \approx_\epsilon \text{Size}(C, 2^n, 2^{\epsilon^{-1}})$ for $\epsilon^{-1} \in \text{Log}$, cf. [17, Lemma 2.14]. If $X \cap t \subseteq 2^{|t|}$ is defined by a circuit C and $\epsilon^{-1} \in \text{Log}$, we can define

$$\Pr_{x < t}[x \in X]_\epsilon := \frac{1}{t} \text{Size}(C, 2^{|t|}, 2^{\epsilon^{-1}}).$$

The presented definitions of approximate counting are well-behaved:

► **Proposition 10** (Jeřábek [17]). *(in PV_1) Let $X, X', Y, Y', Z \subseteq 2^n$ and $W, W' \subseteq 2^m$ be definable sets, and $\epsilon, \delta < 1$. Then*

- i) $X \subseteq Y \Rightarrow X \preceq_0 Y$,
- ii) $X \preceq_\epsilon Y \wedge Y \preceq_\delta Z \Rightarrow X \preceq_{\epsilon+\delta} Z$,
- iii) $X \preceq_\epsilon X' \wedge W \preceq_\delta W' \Rightarrow X \times W \preceq_{\epsilon+\delta+\epsilon\delta} X' \times W'$.
- iv) $X \preceq_\epsilon X' \wedge Y \preceq_\delta Y'$ and X', Y' are separable by a circuit, then $X \cup Y \preceq_{\epsilon+\delta} X' \cup Y'$.

► **Proposition 11** (Jeřábek [17]). *(in APC_1)*

1. Let $X, Y \subseteq 2^n$ be definable by circuits, $s, t, u \leq 2^n$, $\epsilon, \delta, \theta, \gamma < 1, \gamma^{-1} \in \text{Log}$. Then
 - (i) $X \preceq_\gamma Y$ or $Y \preceq_\gamma X$,
 - (ii) $s \preceq_\epsilon X \preceq_\delta t \Rightarrow s < t + (\epsilon + \delta + \gamma)2^n$,
 - (iii) $X \preceq_\epsilon Y \Rightarrow 2^n \setminus Y \preceq_{\epsilon+\gamma} 2^n \setminus X$,
 - (iv) $X \approx_\epsilon s \wedge Y \approx_\delta t \wedge X \cap Y \approx_\theta u \Rightarrow X \cup Y \approx_{\epsilon+\delta+\theta+\gamma} s + t - u$.
2. (Disjoint union) Let $X_i \subseteq 2^n$, $i < m$ be defined by a sequence of circuits and $\epsilon, \delta \leq 1$, $\delta^{-1} \in \text{Log}$. If $X_i \preceq_\epsilon s_i$ for every $i < m$, then $\bigcup_{i < m} (X_i \times \{i\}) \preceq_{\epsilon+\delta} \sum_{i < m} s_i$.
3. (Averaging) Let $X \subseteq 2^n \times 2^m$ and $Y \subseteq 2^m$ be definable by circuits, $Y \preceq_\epsilon t$ and $X_y \preceq_\delta s$ for every $y \in Y$, where $X_y := \{x \mid \langle x, y \rangle \in X\}$. Then for any $\gamma^{-1} \in \text{Log}$,

$$X \cap (2^n \times Y) \preceq_{\epsilon+\delta+\epsilon\delta+\gamma} st.$$

When proving Σ_1^b statements in APC_1 we can afford to work in $\text{S}_2^1 + dWPHP(\text{PV}) + \text{BB}(\Sigma_2^b)$ and, in fact, assuming the existence of a single hard function in PV_1 gives us the full power of APC_1 . Here, $\text{BB}(\Sigma_2^b)$ is defined as $\text{BB}(\Sigma_1^b)$ but with $A \in \Sigma_2^b$.

► **Lemma 12** ([25]). *Suppose $\text{S}_2^1 + dWPHP(\text{PV}) + \text{BB}(\Sigma_2^b) \vdash \exists y A(x, y)$ for $A \in \Sigma_1^b$. Then, for every $\epsilon < 1$, there is k and PV-functions g, h such that PV_1 proves*

$$|f| \geq |x|^k \wedge \exists y, |y| = \|f\|, C_h(y) \neq f(y) \rightarrow A(x, g(x, f))$$

where $f(y)$ is the y th bit of f , $f(y) = 0$ for $y > |f|$, and C_h is a circuit of size $\leq 2^{\epsilon\|f\|}$ generated by h on f, x . Moreover, $\text{APC}_1 \vdash \exists y A(x, y)$.

3 Formalizing complexity-theoretic statements

3.1 Circuit lower bounds

An “almost everywhere” formulation of a circuit lower bound for circuits of size s and a function f says that for every sufficiently big n , for each circuit C with n inputs and size s , there exists an input y on which the circuit C fails to compute $f(y)$.

If $f : \{0, 1\}^n \rightarrow \{0, 1\}$ is an NP function and $s = n^k$ for a constant k , this can be written down as a $\forall\Sigma_2^b$ formula $\text{LB}(f, n^k)$,

$$\forall n, n > n_0 \forall \text{ circuit } C \text{ of size } \leq n^k \exists y, |y| = n, C(y) \neq f(y),$$

where n_0 is a constant and $C(y) \neq f(y)$ is a Σ_2^b formula stating that a circuit C on input y outputs the opposite value of $f(y)$.

If we want to express $s(n)$ -size lower bounds for $s(n)$ as big as $2^{O(n)}$, we add an extra assumption on n stating that $\exists x, n = ||x||$. That is, the resulting formula $\text{LB}_{\text{tt}}(f, s(n))$ has form “ $\forall x, n; n = ||x|| \rightarrow \dots$ ”. Treating x, n as free variables, $\text{LB}_{\text{tt}}(f, s(n))$ is Π_1^b if f is, for instance, SAT because $n = ||x||$ implies that the quantifiers bounded by $2^{O(n)}$ are sharply bounded. Moreover, allowing $f \in \text{NE}$ lifts the complexity of $\text{LB}_{\text{tt}}(f, s(n))$ just to $\forall\Sigma_1^b$. The function $s(n)$ in $\text{LB}_{\text{tt}}(f, s(n))$ is assumed to be a PV-function with input x (satisfying $||x|| = n$).

In terms of the *Log*-notation, $\text{LB}(f, n^k)$ implicitly assumes $n \in \text{Log}$ while $\text{LB}_{\text{tt}}(f, n^k)$ assumes $n \in \text{LogLog}$. By choosing the scale of n we are determining how big objects are going to be “feasible” for theories reasoning about the statement. In the case $n \in \text{LogLog}$, the truth-table of f (and everything polynomial in it) is feasible. Assuming just $n \in \text{Log}$ means that only the objects of polynomial-size in the size of the circuit are feasible. Likewise, the theory reasoning about the circuit lower bound is less powerful when working with $\text{LB}(f, n^k)$ than with $\text{LB}_{\text{tt}}(f, n^k)$. (The scaling in $\text{LB}_{\text{tt}}(f, s)$ corresponds to the choice of parameters in natural proofs and in the formalizations by Razborov [29].)

We can analogously define formulas $\text{LB}_{\text{tt}}(f, s(n), t(n))$ expressing an average-case lower bound for f , where f is a free variable (with $f(y)$ being the y th bit of f and $f(y) = 0$ for $y > |f|$). More precisely, $\text{LB}_{\text{tt}}(f, s(n), t(n))$ generalizes $\text{LB}_{\text{tt}}(f, s(n))$ by saying that each circuit of size $s(n)$ fails to compute f on at least $t(n)$ inputs, for PV-functions $s(n), t(n)$. Since $n \in \text{LogLog}$, $\text{LB}_{\text{tt}}(f, s(n), t(n))$ is Π_1^b .

Propositional version. An $s(n)$ -size circuit lower bound for a function $f : \{0, 1\}^n \rightarrow \{0, 1\}$ can be expressed by a $2^{O(n)}$ -size propositional formula $\text{tt}(f, s)$,

$$\bigvee_{y \in \{0, 1\}^n} f(y) \neq C(y)$$

where the formula $f(y) \neq C(y)$ says that an $s(n)$ -size circuit C represented by $\text{poly}(s)$ variables does not output $f(y)$ on input y . The values $f(y)$ are fixed bits. That is, the whole truth-table of f is hard-wired in $\text{tt}(f, s)$.

The details of the encoding of the formula $\text{tt}(f, s)$ are not important for us as long as the encoding is natural because systems like EF considered in this paper can reason efficiently about them. We will assume that $\text{tt}(f, s)$ is the formula resulting from the translation of Π_1^b formula $\text{LB}_{\text{tt}}(h, s)$, where $n_0 = 0$, n, x are substituted after the translation by fixed constants so that $x = 2^{2^n}$, and h is a free variable (with $h(y)$ being the y th bit of h and $h(y) = 0$ for $y > |h|$) which is substituted after the translation by constants defining f .

Analogously, we can express average-case lower bounds by propositional formulas $\text{tt}(f, s(n), t(n))$ obtained by translating $\text{LB}_{\text{tt}}(h, s(n), t(n)2^n)$, with $n_0 = 0$, fixed $x = 2^{2^n}$ and h substituted after the translation by f .

3.2 Learning algorithms

A circuit class \mathcal{C} is defined by a PV-formula if there is a PV-formula defining the predicate $C \in \mathcal{C}$. Definition 6 can be formulated in the language of HARD^A : A circuit class \mathcal{C} (defined by a PV-formula) is learnable over the uniform distribution by a circuit class \mathcal{D} (defined by a PV-formula) up to error ϵ with confidence δ , if there are randomized oracle circuits L^f from \mathcal{D} such that for every Boolean function $f : \{0, 1\}^n \mapsto \{0, 1\}$ (represented by its truth-table) computable by a circuit from \mathcal{C} , for each $\gamma^{-1} \in \text{Log}$, when given oracle access to f , input 1^n and the internal randomness $w \in \{0, 1\}^*$, L^f outputs the description of a circuit satisfying

$$\Pr_w[L^f(1^n, w) \text{ } (1 - \epsilon)\text{-approximates } f]_\gamma \geq \delta.$$

The inner probability of approximability of f by $L^f(1^n, w)$ is counted exactly. This is possible because f is represented by its truth-table, which implies that $2^n \in \text{Log}$.⁸

Propositional version. In order, to translate the definition of learning algorithms to propositional formulas we need to look more closely at the definition of HARD^A .

PV_1 can be relativized to $\text{PV}_1(\alpha)$. The new function symbol α is then allowed in the inductive clauses for introduction of new function symbols. This means that the language of $\text{PV}_1(\alpha)$, denoted also $\text{PV}(\alpha)$, contains symbols for all p-time oracle algorithms.

► **Proposition 13** (Jeřábek [15]). *For every constant $\epsilon < 1/3$ there exists a constant n_0 such that APC_1 proves: for every $n \in \text{LogLog}$ such that $n > n_0$, there exist a function $f : 2^n \rightarrow 2$ such that no circuit of size $2^{\epsilon n}$ computes f on $\geq (1/2 + 1/2^{\epsilon n})2^n$ inputs.*

► **Definition 14** (Jeřábek [15]). *The theory HARD^A is an extension of the theory $\text{PV}_1(\alpha) + \text{dWPHP}(\text{PV}(\alpha))$ by the axioms*

1. $\alpha(x)$ is a truth-table of a Boolean function in $\|x\|$ variables,
2. $\text{LB}_{\text{tt}}(\alpha(x), 2^{\|x\|/4}, 2^{\|x\|}(1/2 - 1/2^{\|x\|/4}))$, for constant n_0 from Proposition 13,
3. $\|x\| = \|y\| \rightarrow \alpha(x) = \alpha(y)$.

By inspecting the proof of Lemma 2.14 in [17], we can observe that on each input $C, 2^n, 2^{\epsilon^{-1}}$ the $\text{PV}_1(\alpha)$ -function *Size* calls α just once (to get the truth-table of a hard function which is then used as the base function of the Nisan-Wigderson generator). In fact, *Size* calls α on input x which depends only on $|C|$, the number of inputs of C and w.l.o.g. also just on $\|\epsilon^{-1}\|$ (since decreasing ϵ leads only to a better approximation). In combination with the fact that the approximation $\text{Size}(C, 2^n, 2^{\epsilon^{-1}}) \approx_\epsilon X$, for $X \subseteq 2^n$ defined by C , is not affected by a particular choice of the hard boolean function generated by α , we get that APC_1 proves

$$\text{LB}_{\text{tt}}(y, 2^{\|y\|/4}, 2^{\|y\|}(1/2 - 1/2^{\|y\|/4})) \wedge \|y\| = S(C, 2^n, 2^{\epsilon^{-1}}) \rightarrow Sz(C, 2^n, 2^{\epsilon^{-1}}, y) \approx_\epsilon X,$$

where Sz is defined as *Size* with the only difference that the call to $\alpha(x)$ on $C, 2^n, 2^{\epsilon^{-1}}$ is replaced by y and $S(C, 2^n, 2^{\epsilon^{-1}}) = \|x\|$ for a PV-function S .

⁸ It could be interesting to develop systematically a standard theory of learning algorithms in APC_1 and WF , but it is not our goal here. Note, for example, that when we are learning small circuits it is not clear how to boost the confidence to 1 in APC_1 , because we don't have counting with exponential precision.

This allows us to express $\Pr_{x < t}[x \in X]_\epsilon = a$, where $\epsilon^{-1} \in \text{Log}$ and $X \cap t \subseteq 2^{|t|}$ is defined by a circuit C , without a $\text{PV}_1(\alpha)$ function, by formula

$$\forall y (\text{LB}_{\text{tt}}(y, 2^{\lceil |y|/4 \rceil}, 2^{\lceil |y| \rceil} (1/2 - 1/2^{\lceil |y|/4 \rceil})) \wedge \|y\| = S(C, 2^{|t|}, 2^{\epsilon^{-1}}) \rightarrow Sz(C, 2^{|t|}, 2^{\epsilon^{-1}}, y)/t = a).$$

We denote the resulting formula by $\Pr_{x < t}^y[x \in X]_\epsilon = a$. We will use the notation $\Pr_{x < t}^y[x \in X]_\epsilon$ in equations with the intended meaning that the equation holds for the value $Sz(\cdot, \cdot, \cdot, \cdot)/t$ under corresponding assumptions. For example, $t \cdot \Pr_{x < t}^y[x \in X]_\epsilon \preceq_\delta a$ stands for “ $\forall y, \exists v, \exists$ circuit \hat{C} (defining a surjection) which witnesses that $\text{LB}_{\text{tt}}(y, 2^{\lceil |y|/4 \rceil}, 2^{\lceil |y| \rceil} (1/2 - 1/2^{\lceil |y|/4 \rceil})) \wedge \|y\| = S(\hat{C}, 2^{|t|}, 2^{\epsilon^{-1}})$ implies $Sz(\hat{C}, 2^{|t|}, 2^{\epsilon^{-1}}, y) \preceq_\delta a$ ”.

The definition of learning can be now expressed without a $\text{PV}_1(\alpha)$ function: If circuit class \mathcal{C} is defined by a PV-function, the statement that a given oracle algorithm L (given by a PV-function with oracle queries) learns a circuit class \mathcal{C} over the uniform distribution up to error ϵ with confidence δ can be expressed as before with the only difference that we replace $\Pr_w[L^f(1^n, w) (1 - \epsilon)\text{-approximates } f]_\gamma \geq \delta$ by

$$\Pr_w^y[L^f(1^n, w) (1 - \epsilon)\text{-approximates } f]_\gamma \geq \delta.$$

Since the resulting formula A defining learning is not Π_1^b (because of the assumption LB_{tt}) we cannot translate it to propositional logic. We will sidestep the issue by translating only the formula B obtained from A by deleting subformula LB_{tt} (but leaving $\|y\| = S(\cdot, \cdot, \cdot)$ intact) and replacing the variables y by fixed bits representing a hard boolean function. In more detail, Π_1^b formula B can be translated into a sequence of propositional formulas $\text{lear}_\gamma^y(L, \mathcal{C}, \epsilon, \delta)$ expressing that “if $C \in \mathcal{C}$ is a circuit computing f , then L querying f generates a circuit D such that $\Pr[D(x) = f(x)] \geq 1 - \epsilon$ with probability $\geq \delta$, which is counted approximately with precision γ ”. Note that C, f are represented by free variables and that there are also free variables for error γ from approximate counting and for boolean functions y . As in the case of tt -formulas, we fix $|f| = 2^n$, so n is not a free variable. Importantly, $\text{lear}_\gamma^y(L, \mathcal{C}, \epsilon, \delta)$ does not postulate that y is a truth-table of a hard boolean function. Nevertheless, for any fixed (possibly non-uniform) bits representing a sequence of boolean functions $h = \{h_m\}_{m > n_0}$ such that h_m is not $(1/2 + 1/2^{m/4})$ -approximable by any circuit of size $2^{m/4}$, we can obtain formulas $\text{lear}_\gamma^h(L, \mathcal{C}, \epsilon, \delta)$ by substituting bits h for y .

Using a single function h in $\text{lear}_\gamma^h(L, \mathcal{C}, \epsilon, \delta)$ does not ruin the fact that (the translation of function) Sz approximates the respective probability with accuracy γ because Sz queries a boolean function y which depends just on the number of atoms representing γ^{-1} and on the size of the circuit D defining the predicate we count together with the number of inputs of D . The size of D and the number of its inputs are w.l.o.g. determined by the number of inputs of f .

If we are working with formulas $\text{lear}_\gamma^h(L, \mathcal{C}, \epsilon, \delta)$, where h is a sequence of bits representing a hard boolean function, in a proof system which cannot prove efficiently that h is hard, our proof system might not be able to show that the definition is well-behaved - it might not be able to derive some standard properties of the function Sz used inside the formula. Nevertheless, in our theorems this will never be the case: our proof systems will always know that h is hard.

In formulas $\text{lear}_\gamma^y(L, \mathcal{C}, \epsilon, \delta)$ we can allow L to be a sequence of nonuniform circuits, with a different advice string for each input length. One way to see that is to use additional input to L in Π_1^b formula B , then translate the formula to propositional logic and substitute the right bits of advice for the additional input. Again, the precise encoding of the formula $\text{lear}_\gamma^y(L, \mathcal{C}, \epsilon, \delta)$ does not matter very much to us but in order to simplify proofs we will

assume that $\text{lear}_\gamma^y(L, \text{Circuit}[n^k], \epsilon, \delta)$ has the form $\neg \text{tt}(f, n^k) \rightarrow R$, where n, k are fixed, f is represented by free variables and R is the remaining part of the formula expressing that L generates a suitable circuit with high probability.

3.3 Automatability

Let Φ be a class of propositional formulas. We say that a proof system P is automatable w.r.t. Φ up to proofs of size s , where $s : \Phi \mapsto \mathbb{N}$ is a function, if there is a PV-function A such that for each $\phi \in \Phi$ and each t -size P -proof of ϕ with $t \leq s$, $A(\phi, 1^t)$ is a P -proof of ϕ .

In our main theorem we will need a slightly modified notion of automatability where the automating algorithm outputs a proof of a given tautology ϕ which is not much longer than a proof of an associated tautology ψ . (Formula ψ will be closely related to ϕ : while ϕ will express a worst-case lower bound, ψ will express an average-case lower bound for the same function.)

Let Φ be a class of pairs of propositional formulas. We say that a proof system P is automatable w.r.t. Φ up to proofs of size s if there is a PV-function A such that for each pair $\langle \psi, \phi \rangle \in \Phi$ and each t -size P -proof of ψ with $t \leq s$, $A(\phi, 1^t)$ is a P -proof of ϕ .

Propositional version. If Φ is defined by a PV-function and $s \in \text{Log}$ is a PV-function, the statement that an algorithm A (given by a PV-function) automates system P w.r.t. Φ up to proofs of size s is Π_1^b . Therefore, it can be translated into a sequence of propositional formulas $\text{aut}_P(A, \Phi, s)$. Again, in formulas $\text{aut}_P(A, \Phi, s)$ we can allow A to be a sequence of nonuniform circuits and s to be arbitrary possibly nonuniform parameter.

4 Learning algorithms from natural proofs in APC_1

The formalization of the transformation of natural proofs into learning algorithms follows from a straightforward inspection of the original proof. The proof can be found in the arXiv version of the paper.

► **Theorem 15.** *There is a PV-function L such that APC_1 proves: For $k \geq 1$, $d \geq 2$, $2^{n^d}, n^{dk}, \delta^{-1} \in \text{Log}$, $\delta < 1/N^3$ and a prime $n^d \leq p \leq 2n^d$, let R_N be a circuit with $N = 2^n$ inputs such that for sufficiently big N ,*

1. $R_N(x) = 1$ implies that x is a truth-table of a boolean function with n inputs hard for $\text{Circuit}[n^{10dk}]$,
2. $\{x \mid R_N(x) = 1\} \succeq_\delta 2^N/N$.

Then, circuits with n^d inputs and size n^{dk} are learnable by circuit $L(R_N, p)$ over the uniform distribution with membership queries, confidence $1/N^4$, up to error $1/2 - 1/N^3$. Here, the confidence is counted approximately with error δ using PV-function S_z and the corresponding assumptions LB_{tt} expressing hardness of a boolean function y , i.e. using formulas $\text{Pr}^y[\cdot]_\delta$.

5 Main theorem

Our main theorem holds for any “decent” proof system p -simulating WF, which is well-behaved in the sense that it APC_1 -provably satisfies some basic properties.

► **Definition 16** (APC_1 -decent proof system). *A propositional proof system P is APC_1 -decent if the language L of P is finite and complete, i.e. L consists of connectives of constant arity such that each boolean function of every arity can be expressed by an L -formula, P proves efficiently its own reflection principle, i.e. formulas stating that if π is a P -proof of ϕ then ϕ holds, cf. [22], and there is a PV-function F such that APC_1 proves:*

1. P p -simulates WF, i.e. F maps each WF-proof of ϕ to a P -proof of ϕ .
2. P admits substitution property: F maps each triple $\langle \phi, \rho, \pi \rangle$ to a P -proof of $\phi|_\rho$, where π is a P -proof of ϕ and $\phi|_\rho$ is formula ϕ after applying substitution ρ which replaces atoms of ϕ by formulas.
3. F maps each pair $\langle \pi, \pi' \rangle$, where π is a P -proof of ϕ and π' is a P -proof of $\phi \rightarrow \psi$, to a P -proof of ψ .

In Definition 16, WF refers to some fixed system from the set of all WF systems. It follows from the proof of Lemma 17 that if APC_1 proves that P p -simulates a WF-system Q , then for every WF-system R , APC_1 proves that P p -simulates R , so the particular choice of the WF-system does not matter. When we use connectives $\wedge, \vee, \neg, \rightarrow$ in an APC_1 -decent system P , we assume that these are expressed in the language of P .

► **Lemma 17.** *Each WF system is APC_1 -decent. Moreover, for each APC_1 -decent proof system P the following holds.*

1. For every Frege rule which derives ϕ from ϕ_1, \dots, ϕ_k , there is a PV-function F such that APC_1 proves that F maps each $(k+1)$ -tuple $\langle \pi_1, \dots, \pi_k, \rho \rangle$ to a P -proof of $\phi|_\rho$, where π_i is a P -proof of $\phi_i|_\rho$ for a substitution ρ replacing each atom of $\phi, \phi_1, \dots, \phi_k$ by a formula.
2. There is a PV-function F such that APC_1 proves that F maps each pair $\langle \phi, b \rangle$, for assignment b satisfying formula ϕ , to a P -proof of $\phi(b)$.
3. Let π be a P -proof of $E \rightarrow \phi$, where E defines a computation of a circuit which is allowed to use atoms from ϕ as inputs but other atoms of E do not appear in ϕ , i.e. E is the conjunction of extension axioms of EF built on atoms from ϕ . Then, there is a $\text{poly}(|\pi|)$ -size P -proof of ϕ .

Proof. WF is known to prove efficiently its own reflection principle, cf. [15]. In order to show that it is APC_1 -decent, it thus suffices to prove that it satisfies Items 1-3 from Definition 16.

Item 2 is established already in PV_1 by Σ_1^b -induction on the length of the proof π (which can be used because of $\forall \Sigma_1^b$ -conservativity of S_2^1 over PV_1): F replaces each circuit C from π by $C|_\rho$ and preserves all WF-derivation rules.

Item 1 holds trivially if the given WF-system P is the WF-system P' from Definition 16. Otherwise, we use implicational completeness of P and the completeness of the language of P to simulate all $O(1)$ Frege rules of P' by $O(1)$ steps in P . (This does not require that the implicational completeness of P is provable in APC_1 because we need to simulate only $O(1)$ Frege rules of finite size). Similarly, by Σ_1^b -induction and the completeness of the language of P , we simulate each circuit in the language of P' by a circuit in the language of P and show that this simulation preserves the similarity rule. Then, given an s -size P' -proof of ϕ , we obtain a $\text{poly}(s)$ -size P -proof of ϕ using the simulation of Frege rules of P' , the similarity rule and dWPHP axiom, together with substituting the right circuits in Frege rules. This is done again in PV_1 by Σ_1^b -induction on the length of the P' -proof.

Item 3 follows by simulating modus ponens as in the proof of Item 1.

For the “moreover” part, see the arXiv version of the paper. ◀

APC_1 -decent proof systems can be much stronger than WF. For example, consider ZFC as a propositional proof system: a ZFC-proof of propositional formula ϕ is a ZFC-proof of the statement encoding that ϕ is a tautology. We can add the reflection of ZFC to WF, i.e. we will allow WF to derive (substitutional instances of) formulas stating that “If π is a ZFC-proof of ϕ , then ϕ holds.” The new system is as strong as ZFC w.r.t. tautologies and it is easy to see that it is APC_1 -decent. (The reflection of the system can be proved in APC_1 extended with an axiom postulating the reflection for ZFC.)

► **Theorem 18** (Learning versus automatability). *Let P be an APC_1 -decent proof system and assume there is a sequence of boolean functions $h = \{h_n\}_{n > n_1}$, for a constant n_1 , such that P proves efficiently $\text{tt}(h_n, 2^{n/4}, 1/2 - 1/2^{n/4})$. Then, for each constant K and constant $\gamma < 1$, the following statements are equivalent.*

1. **Provable learning.** *For each $k \geq 1$ and $\ell \geq K + 1$, there are 2^{Kn^γ} -size circuits A such that for each sufficiently big n , P proves efficiently*

$$\text{lear}_{1/2^{\ell n^\gamma}}^h(A, \text{Circuit}[n^k], 1/2 - 1/2^{Kn^\gamma}, 1/2^{Kn^\gamma}).$$

2. **Provable automatability.** *For each $k \geq 1$, for each function $s(n) \geq 2^n$, there is a constant K' and $s^{K'}$ -size circuits B such that P proves efficiently*

$$\text{aut}_P(B, \Phi, s),$$

where Φ is the set of pairs $\langle \text{tt}(f, 2^{Kn^\gamma}, 1/2 - 1/2^{Kn^\gamma}), \text{tt}(f, n^k) \rangle$ for all boolean functions f with n inputs.

Proof. (1. \rightarrow 2.) We first prove the following statement in APC_1 .

▷ **Claim 19** (in APC_1). Assume that π is a P -proof of $\text{lear}_{1/2^{\ell n^\gamma}}^y(A, \text{Circuit}[n^k], 1/2 - 1/2^{Kn^\gamma}, \delta)$ for a circuit A and a boolean function y represented by fixed bits in formula $\text{lear}_{1/2^{\ell n^\gamma}}^y(\cdot, \cdot, \cdot, \cdot)$. Further, assume that the probability that A on queries to f outputs a circuit D such that $\Pr[D(x) = f(x)] \geq 1/2 + 1/2^{Kn^\gamma}$ is $< \delta$, where the outermost probability is counted approximately with error $1/2^{\ell n^\gamma}$ using PV-function Sz and the corresponding assumptions LB_{tt} expressing hardness of y , i.e. using formulas $\text{Pr}^y[\cdot]_{1/2^{\ell n^\gamma}}$ for the same y as above - we treat y as a free variable here. Then there is a $\text{poly}(|\pi|)$ -size P -proof of $\text{tt}(f, n^k)$ or y does not satisfy the assumptions of $\text{Pr}^y[\cdot]_{1/2^{\ell n^\gamma}}$.

To see that the claim holds, we reason in APC_1 as follows. Assume π is a P -proof of $\text{lear}_{1/2^{\ell n^\gamma}}^y(A, \text{Circuit}[n^k], 1/2 - 1/2^{Kn^\gamma}, \delta)$ but A on queries to f outputs a circuit $(1/2 + 1/2^{Kn^\gamma})$ -approximating f with probability $< \delta$. Then, either y does not satisfy the assumptions of $\text{Pr}^y[\cdot]_{1/2^{\ell n^\gamma}}$ or there is a trivial $2^{O(n)}$ -size P -proof of $\neg \text{tt}(f, n^k) \rightarrow \neg R(b)$, for predicate R from the definition of $\text{lear}_{1/2^{\ell n^\gamma}}^y(A, \text{Circuit}[n^k], 1/2 - 1/2^{Kn^\gamma}, \delta)$ and a complete assignment b . The P -proof is obtained by evaluating function Sz which counts the confidence of A - note that functions f, y and algorithm A are represented inside P by fixed bits so the P -proof just evaluates a $2^{O(n)}$ -size circuit on some input, which is possible by Lemma 17, Item 2. (We use here also the fact that APC_1 knows that the probability statement expressed by function Sz translates to $\neg R$ in the negation normal form.) The formula $\neg \text{tt}(f, n^k) \rightarrow \neg R(b)$ is obtained from $\neg R(b)$ by an instantiation of a single Frege rule, which is available by Lemma 17, Item 1. Applying again Lemma 17, Item 1, from a P -proof of $\text{lear}_{1/2^{\ell n^\gamma}}^y(A, \text{Circuit}[n^k], 1/2 - 1/2^{Kn^\gamma}, \delta)$ and a P -proof of $\neg \text{tt}(f, n^k) \rightarrow \neg R(b)$, we construct a $\text{poly}(|\pi|)$ -size P -proof of $\text{tt}(f, n^k)$. This proves the claim.

Next, observe that APC_1 proves that “If for a sufficiently big n and $\ell \geq K + 1$ the probability that a circuit A on queries to f outputs a circuit $(1/2 + 1/2^{Kn^\gamma})$ -approximating f is $\geq 1/2^{Kn^\gamma}$, where the probability is counted approximately with error $1/2^{\ell n^\gamma}$ using PV-function Sz and the corresponding assumptions LB_{tt} , then there is a circuit of size $|A|$ $(1/2 + 1/2^{Kn^\gamma})$ -approximating f or y does not satisfy the assumptions of $\text{Pr}^y[\cdot]_{1/2^{\ell n^\gamma}}$.” This is because, if such a circuit did not exist, a trivial surjection would witness that 2^m times the probability that A outputs a circuit $(1/2 + 1/2^{Kn^\gamma})$ -approximating f , counted approximately with error $1/2^{\ell n^\gamma}$ using function Sz , is $\leq_{1/2^{\ell n^\gamma}} 0$. Here, 2^m is the domain of the surjection. By Proposition 11 1.ii), this would imply $2^m/2^{Kn^\gamma} < 2^{m+1}/2^{\ell n^\gamma}$, which is a contradiction for $\ell \geq K + 1$ and sufficiently big n .

Therefore, Claim 19 implies that APC_1 proves that “For sufficiently big n and $\ell \geq K + 1$, if π is a P -proof of $\text{lear}_{1/2^{\ell n^\gamma}}^y(A, \text{Circuit}[n^k], 1/2 - 1/2^{Kn^\gamma}, 1/2^{Kn^\gamma})$ for circuits A of size 2^{Kn^γ} , then there is a P -proof of $\text{tt}(f, n^k)$ or there is a 2^{Kn^γ} -size circuit $(1/2 + 1/2^{Kn^\gamma})$ -approximating f or there is a $2^{\|y\|/4}$ -size circuit $(1/2 + 1/2^{\|y\|/4})$ -approximating y or $\|y\| \leq n_0$ or $\|y\| \neq S(\cdot, 2^m, 2^{2^{\ell n^\gamma}})$,” for n_0 from Definition 14. Since this is a Σ_1^b -statement, by Lemma 12, PV_1 proves the same statement with the existential quantifiers witnessed by PV -functions assuming they are given a boolean function h' which is hard for circuits of size $2^{\|h'\|/4}$, for sufficiently big $|h'|$.

The last statement provable in PV_1 is Π_1^b so we can translate it to EF . This gives us $\text{poly}(|\pi|, 2^n)$ -size circuits B_0 such that for sufficiently big n , EF proves efficiently

“If $\ell \geq K + 1$,

h' is not computable by a particular circuit of size $2^{\|h'\|/4}$, $|h'|$ is sufficiently big,

y is not $(1/2 + 1/2^{\|y\|/4})$ -approximable by a particular circuit of size $2^{\|y\|/4}$, $\|y\| > n_0$,

$\|y\| = S(\cdot, 2^m, 2^{2^{\ell n^\gamma}})$

and π is a P -proof of $\text{lear}_{1/2^{\ell n^\gamma}}^y(A, \text{Circuit}[n^k], 1/2 - 1/2^{Kn^\gamma}, 1/2^{Kn^\gamma})$ for 2^{Kn^γ} -size A ,

then B_0 (given π, h' and formula $\text{tt}(f, n^k)$) outputs a P -proof of $\text{tt}(f, n^k)$

or B_0 outputs a 2^{Kn^γ} -size circuit $(1/2 + 1/2^{Kn^\gamma})$ -approximating f .”⁹

If we now assume that P proves efficiently $\text{tt}(h_n, 2^{n/4}, 1/2 - 1/2^{n/4})$ and that Item 1 holds, then by Definition 16, Items 1-3, for each k , there are p -size circuits B_1 such that for each sufficiently big n , P proves efficiently “ B_1 (given just formula $\text{tt}(f, n^k)$) outputs a P -proof of $\text{tt}(f, n^k)$ or B_1 outputs a 2^{Kn^γ} -size circuit $(1/2 + 1/2^{Kn^\gamma})$ -approximating f .” (We use here also the fact that PV_1 knows that $S(\cdot, 2^m, 2^{2^{\ell n^\gamma}})$ depends just on n .) Consequently, since P proves efficiently its own reflection, for each sufficiently big n , P proves efficiently that “if π is a P -proof of $\text{tt}(f, 2^{Kn^\gamma}, 1/2 - 1/2^{Kn^\gamma})$ then B_1 outputs a P -proof of $\text{tt}(f, n^k)$.”¹⁰ Finally, we make the P -proofs work for all n by increasing the size of B_1 by a constant. This finishes the proof of case (1. \rightarrow 2.).

(2. \rightarrow 1.) The opposite implication can be obtained from Lemma 20 and 21 which formalize Theorem 4.

► **Lemma 20.** *For each $d \geq 2$, each $k \geq 10d$ and each sufficiently big c , there is a PV -function L such that for each PV -function B the theory APC_1 proves: Assume the reflection principle for P holds, π is a P -proof of*

$$\text{tt}(h_n \oplus g, 2^{Kn^\gamma}, 1/2 - 1/2^{Kn^\gamma}) \vee \text{tt}(g, 2^{Kn^\gamma}), \quad (1)$$

where g is represented by free variables, and that B automates P on Φ up to size $|\pi|^c$. Then, for prime $n^d \leq p \leq 2n^d$, where $2^{n^d} \in \text{Log}$, for $\delta^{-1} \in \text{Log}$ such that $\delta < 1/N^3 = 2^{3n}$, $L(B, \pi, p)$ is a $\text{poly}(2^n, |\pi|)$ -size circuit learning circuits with $m = n^d$ inputs and size $m^{k/10d}$,

⁹ Formally, the statement “If a particular assignment a satisfies formula ϕ , then formula ψ holds” means that “If a is the output of a computation of a specific circuit W (where W is allowed to use as inputs atoms from ψ , but other atoms of W do not appear in ψ), and a satisfies ϕ , then ψ ”. By Lemma 17, Item 3, if we assume that the statement is efficiently provable in P and that P proves efficiently ϕ , then P proves efficiently ψ . Note also that for $A, B \in \Sigma_0^b$, the translation $\|A \rightarrow B\|$ is $\neg\|\neg A\| \rightarrow \|B\|$, which might not be the same formula as $\|A\| \rightarrow \|B\|$. Nevertheless, EF proves efficiently that $E \rightarrow (\|A\| \leftrightarrow \neg\|\neg A\|)$, where E postulates that auxiliary variables of $\|A\|$ encode the computation of a suitable circuit. Therefore, in systems like EF or P , if we have a proof of $\|A\|$ and $\|A \rightarrow B\|$, we can remove the assumption E after proving $E \rightarrow \|B\|$, assuming “non-input” variables of E do not occur in $\|B\|$, and ignore the difference between $\|A\|$ and $\neg\|\neg A\|$.

¹⁰ It is assumed that the encoding of the statement coincides with the encoding of aut_P .

with confidence $1/N^4$, up to error $1/2 - 1/N^3$, where the confidence is counted approximately with error δ using PV-function Sz and the corresponding assumptions LB_{tt} expressing hardness of a boolean function y , i.e. using formulas $\text{Pr}^y[\cdot]_\delta$.

► **Lemma 21** (“XOR trick”). PV_1 proves that for all boolean functions g, h'' with n inputs, for sufficiently big n , $\text{LB}_{\text{tt}}'(h'', 3 \cdot 2^{Kn^\gamma}, 2^n(1/2 - 1/2^{Kn^\gamma}))$ implies $\text{LB}_{\text{tt}}'(h'' \oplus g, 2^{Kn^\gamma}, 2^n(1/2 - 1/2^{Kn^\gamma})) \vee \text{LB}_{\text{tt}}'(g, 2^{Kn^\gamma})$, where LB_{tt}' is obtained from LB_{tt} by setting $n_0 = 0$ and skipping the universal quantifier on n , i.e. all formulas LB_{tt}' refer to the same n .

The proof of Lemma 21 is almost immediate: By Σ_1^b -induction, a 2^{Kn^γ} -size circuit C_1 computing g and a 2^{Kn^γ} -size circuit C_2 $(1/2 + 1/2^{Kn^\gamma})$ -approximating $h'' \oplus g$ can be combined into a circuit $C_1 \oplus C_2$ of size $3 \cdot 2^{Kn^\gamma}$ which $(1/2 + 1/2^{Kn^\gamma})$ -approximates h'' .

The implication $(2. \rightarrow 1.)$ can be derived from Lemma 20 and 21 as follows. Since the APC_1 -provable statement from Lemma 20 is Σ_1^b , similarly as above, we can witness it and translate to EF at the expense of introducing an additional assumption about the hardness of a boolean function h' . That is, for each p -size circuit B there are $\text{poly}(|\pi|, 2^{n^d})$ -size circuits A and $\text{poly}(|\pi|, 2^{n^d})$ -size EF-proofs of

“If the reflection principle for P is satisfied by a particular assignment,
 π is a P -proof of (1),
 h' is not computable by a particular circuit of size $2^{\|h'\|/4}$, $|h'|$ is sufficiently big,
 y is not $(1/2 + 1/2^{\|y\|/4})$ -approximable by a particular circuit of size $2^{\|y\|/4}$, $\|y\| > n_0$,
 $\|y\| = S(\cdot, \cdot, 2^{\delta^{-1}})$
and $n^d \leq p \leq 2n^d$ is a prime,
then, for $\delta < 1/N^3$, $\text{lear}_\delta^y(L(B, \pi, p), \text{Circuit}(m^{k/10d}), 1/2 - 1/N^3, 1/N^4)$
or $A(B, \pi, h')$ outputs a falsifying assignment of $\text{aut}_P(B, \Phi, |\pi|^c)$.”

Analogously, PV_1 -proof from Lemma 21 yields p -size EF-proofs of the implication “ $\text{tt}(h_n, 3 \cdot 2^{Kn^\gamma}, 1/2 - 1/2^{Kn^\gamma})$ is falsified by a particular assignment or (1) holds”. By the assumption of the theorem, there are p -size P -proofs of $\text{tt}(h_n, 3 \cdot 2^{Kn^\gamma}, 1/2 - 1/2^{Kn^\gamma})$ for sufficiently big n . Hence, by Definition 16, Items 1-3, there are p -size P -proofs of (1) for sufficiently big n . As P proves efficiently also its own reflection, this yields $\text{poly}(2^{n^d})$ -size P -proofs of

“If h' is not computable by a particular circuit of size $2^{\|h'\|/4}$, $|h'|$ is sufficiently big,
 y is not $(1/2 + 1/2^{\|y\|/4})$ -approximable by a particular circuit of size $2^{\|y\|/4}$, $\|y\| > n_0$,
 $\|y\| = S(\cdot, \cdot, 2^{\delta^{-1}})$
and $n^d \leq p \leq 2n^d$ is a prime,
then, for $\delta < 1/N^3$, $\text{lear}_\delta^y(L(B, \pi, p), \text{Circuit}(m^{k/10d}), 1/2 - 1/N^3, 1/N^4)$
or $A(B, \pi, h')$ outputs a falsifying assignment of $\text{aut}_P(B, \Phi, |\pi|^c)$.”

By Bertrand’s postulate there is a prime $n^d \leq p \leq 2n^d$, so EF proves that p is a prime by a trivial $2^{O(n^d)}$ -size proof which verifies all possible divisors. Therefore, choosing $d > 1/\gamma$, Item 2 and p -size P -proofs of $\text{tt}(h_n, 2^{n/4}, 1/2 - 1/2^{n/4})$ imply Item 1.

It remains to prove Lemma 20.

Suppose π is a P -proof of (1). Assuming that B automates P on Φ , we want to obtain a P/poly -natural property useful against $\text{Circuit}[n^k]$. To do so, observe (first, without formalizing it in APC_1) that for each g , B can be used to find a proof of $\text{tt}(h_n \oplus g, n^k)$ or to recognize that $\text{tt}(g, 2^{Kn^\gamma})$ holds - if $\text{tt}(g, 2^{Kn^\gamma})$ was falsifiable, there would exist a $\text{poly}(|\pi|)$ -size P -proof of $\text{tt}(h_n \oplus g, 2^{Kn^\gamma}, 1/2 - 1/2^{Kn^\gamma})$ obtained by substituting the falsifying assignment to the proof of (1) and thus B would find a short proof of $\text{tt}(h_n \oplus g, n^k)$, for

sufficiently big c . Since for random g , both $h_n \oplus g$ and g are random functions, we know that with probability $\geq 1/2$ B finds a proof of $\text{tt}(h_n \oplus g, n^k)$ or with probability $\geq 1/2$ it recognizes that $\text{tt}(g, 2^{Kn^\gamma})$ holds. In both cases, B yields a P/poly-natural property useful against $\text{Circuit}[n^k]$.

Let us formalize reasoning from the previous paragraph in APC_1 . Let $N = 2^n$ and B' be the algorithm which uses B to search for P -proofs of $\text{tt}(h_n \oplus g, n^k)$ or to recognize that $\text{tt}(g, 2^{Kn^\gamma})$ holds. B' uses π to know how long it needs to run B . Assume for the sake of contradiction that

$$G_0 := \{g \oplus h_n \mid B'(g) \text{ outputs a } P\text{-proof of } \text{tt}(h_n \oplus g, n^k)\} \preceq_0 2^N/3$$

$$G_1 := \{g \mid B'(g) \text{ recognizes that } \text{tt}(g, 2^{Kn^\gamma}) \text{ holds}\} \preceq_0 2^N/3.$$

It is easy to construct a surjection S witnessing that $2^N \preceq_0 G_0 \cup G_1$: S maps $g \in G_1$ to g and $g \in G_0$ to $g \oplus h_n$. Following the argument above we conclude that S is a surjection: for each g , either $g \in G_1$ (and $S(g) = g$) or $g \oplus h_n \in G_0$ (and $S(g \oplus h_n) = g$). Here, we use the assumption that APC_1 knows that P admits the substitution property and simulates Frege rules. Thus, by Proposition 10 iv), $2^N \preceq_0 2 \cdot 2^N/3$, which yields a contradiction by Proposition 11 1.ii). Consequently, by Proposition 11 1.i), $G_0 \succeq_\delta 2^N/3$ or $G_1 \succeq_\delta 2^N/3$ for $\delta^{-1} \in \text{Log}$. Since $g \in G_0$ and $g \in G_1$ are decidable by p-size circuits and we assume the reflection principle for P (which implies that G_0 is useful), this means that either G_0 or G_1 defines a P/poly-natural property useful against $\text{Circuit}[n^k]$.

Finally, by the APC_1 -formalization of [8], Theorem 15, we obtain $\text{poly}(2^n, |\pi|)$ -size circuit $L(B, \pi, p)$ learning circuits with $m = n^d$ inputs and size $n^{k/10}$, over the uniform distribution, with membership queries, confidence $1/N^4$, up to error $1/2 - 1/N^3$. ◀

► **Corollary 22.** *Assume there is a $\text{NE} \cap \text{coNE}$ -function $h_n : \{0, 1\}^n \mapsto \{0, 1\}$ such that for each sufficiently big n , h_n is not $(1/2 + 1/2^{n/4})$ -approximable by $2^{n/4}$ -size circuits. Then there is a proof system P (which can be described explicitly¹¹ given the definition of h_n) such that for each constant K and $\gamma < 1$, Items 1 and 2 from Theorem 18 are equivalent. Moreover, the equivalence holds for each APC_1 -decent system simulating P .*

Proof. See the arXiv version of the paper. ◀

References

- 1 M. Alekhnovich, M. Braverman, V. Feldman, A. R. Klivans, and T. Pitassi. Learnability and automatizability. *IEEE Symposium on Foundations of Computer Science*, 2004.
- 2 B. Applebaum, B. Barak, and D. Xiao. On basing lower bounds for learning on worst-case assumptions. *IEEE Symposium on Foundations of Computer Science*, 2008.
- 3 A. Atserias and M. Müller. Automating resolution is NP-hard. *IEEE Symposium on Foundations of Computer Science*, 2019.
- 4 E. Binnendyk, M. Carmosino, A. Kolokolova, R. Ramyaa, and M. Sabin. Learning with distributional inverters. *Algorithmic Learning Theory*, 2022.
- 5 M.L. Bonet, C. Domingo, R. Gavalda, A. Maciel, and T. Pitassi. Non-automatizability of bounded-depth Frege proofs. *Computational Complexity*, 13(1–2):47–68, 2004.
- 6 M.L. Bonet, T. Pitassi, and R. Raz. On interpolation and automatization for Frege proof systems. *SIAM Journal of Computing*, 29(6):1939–1967, 2000.
- 7 S. Buss. *Bounded arithmetic*. Bibliopolis, 1986.

¹¹ More formally, there is a p-time algorithm R such that given predicates H_0, H_1 defining h_n (see the proof of Corollary 22), R outputs a p-time algorithm defining system P .

- 8 M. Carmosino, R. Impagliazzo, V. Kabanets, and A. Kolokolova. Learning algorithms from natural proofs. *IEEE Symposium on Computational Complexity*, 2016.
- 9 A. Cobham. The intrinsic computational difficulty of functions. *Proceedings of the 2nd International Congress of Logic, Methodology and Philosophy of Science*, 1965.
- 10 S.A. Cook. Feasibly constructive proofs and the propositional calculus. *ACM Symposium on Theory of Computing*, 1975.
- 11 S.A. Cook and N. Thapen. The strength of replacement in weak arithmetic. *ACM Transactions on Computational Logic*, 7(4):749–764, 2006.
- 12 S. de Rezende, M. Göös, J. Nördstrom, T. Pitassi, R. Robere, and D. Sokolov. Automating algebraic proof systems is NP-hard. *IEEE Symposium on Computational Complexity*, 2020.
- 13 M. Garlík. Failure of feasible disjunction property for k -DNF resolution and NP-hardness of automating it. *Electronic Colloquium on Computational Complexity*, 2020.
- 14 M. Göös, S. Korothe, I. Mertz, and T. Pitassi. Automating cutting planes is NP-hard. *ACM Symposium on Theory of Computing*, 2020.
- 15 E. Jeřábek. Dual weak pigeonhole principle, Boolean complexity and derandomization. *Annals of Pure and Applied Logic*, 129:1–37, 2004.
- 16 E. Jeřábek. Weak pigeonhole principle and randomized computation. *Ph.D. thesis, Charles University in Prague*, 2005.
- 17 E. Jeřábek. Approximate counting in bounded arithmetic. *Journal of Symbolic Logic*, 72:959–993, 2007.
- 18 J. Krajíček. *Bounded arithmetic, propositional logic, and complexity theory*. Cambridge University Press, 1995.
- 19 J. Krajíček. Interpolation theorems, lower bounds for proof systems, and independence results for bounded arithmetic. *Journal of Symbolic Logic*, 66(2):457–486, 1997.
- 20 J. Krajíček. On the weak pigeonhole principle. *Fundamenta Mathematicae*, 170(1–3):123–140, 2001.
- 21 J. Krajíček. *Forcing with random variables and proof complexity*. Cambridge University Press, 2011.
- 22 J. Krajíček. *Proof complexity*. Cambridge University Press, 2019.
- 23 J. Krajíček and P. Pudlák. Some consequences of cryptographical conjectures for S_2^1 and EF. *Information and Computation*, 140(1):82–94, 1998.
- 24 J. Krajíček, P. Pudlák, and G. Takeuti. Bounded arithmetic and the polynomial hierarchy. *Annals of Pure and Applied Logic*, 52:143–153, 1991.
- 25 M. Müller and J. Pich. Feasibly constructive proofs of succinct weak circuit lower bounds. *Annals of Pure and Applied Logic*, 2019.
- 26 I.C. Oliveira and R. Santhanam. Conspiracies between learning algorithms, circuit lower bounds, and pseudorandomness. *IEEE Symposium on Computational Complexity*, 2017.
- 27 J. Pich. Learning algorithms from circuit lower bounds. *preprint*, 2020.
- 28 A.A. Razborov. On provably disjoint NP pairs. *BRICS*, 1994.
- 29 A.A. Razborov. Bounded arithmetic and lower bounds in boolean complexity. *Feasible Mathematics II*, pages 344–386, 1995.
- 30 A.A. Razborov. Unprovability of lower bounds on the circuit size in certain fragments of bounded arithmetic. *Izvestiya of the Russian Academy of Science*, 59:201–224, 1995.
- 31 A.A. Razborov. Pseudorandom generators hard for k -DNF Resolution and Polynomial Calculus. *Annals of Mathematics*, 181(2):415–472, 2015.
- 32 A.A. Razborov and S. Rudich. Natural proofs. *Journal of Computer and System Sciences*, 55(1):24–35, 1997.
- 33 R. Santhanam. Pseudorandomness and the Minimum Circuit Size Problem. *Innovations in Theoretical Computer Science*, 2020.
- 34 L. Valiant. A theory of the learnable. *Communications of the ACM*, 27, 1984.
- 35 R. Williams. Non-uniform ACC circuit lower bounds. *IEEE Symposium on Computational Complexity*, 2011.