

# A remark on pseudo proof systems and hard instances of the satisfiability problem

Jan Maly<sup>1\*</sup> and Moritz Müller<sup>2,3</sup>

<sup>1</sup> Institute of Logic and Computation, Technische Universität Wien, Favoritenstraße 9, 1040 Wien, Austria

<sup>2</sup> Kurt Gödel Research Center, University of Vienna, Währinger Straße 25, 1090 Wien, Austria

<sup>3</sup> Computer Science Department, Universitat Politècnica de Catalunya, Omega-327, Campus Nord, c/ Jordi Girona 1-3, 08034 Barcelona, Spain

Received 3 March 2017, revised 8 September 2017, accepted 17 January 2018

Published online 20 November 2018

We link two concepts from the literature, namely hard sequences for the satisfiability problem SAT and so-called pseudo proof systems proposed for study by Krajíček. Pseudo proof systems are elements of a particular nonstandard model constructed by forcing with random variables. We show that the existence of *mad* pseudo proof systems is equivalent to the existence of a randomized polynomial time procedure with a highly restrictive use of randomness which produces satisfiable formulas whose satisfying assignments are probably hard to find.

© 2018 The Authors. *Mathematical Logic Quarterly* published by WILEY-VCH Verlag GmbH & Co. KGaA, Weinheim

## 1 Introduction

**Pseudo proof systems.** It is a basic question of mathematical logic, unsettled to date, whether there exists a propositional proof system that has *short* proofs for all tautologies. Abstractly, a *propositional proof system* is a polynomial time function from the set of binary strings  $\{0, 1\}^*$  into the set TAUT of (binary strings coding) propositional tautologies. Often [11] it is additionally required, that the function is not only into TAUT, meaning *soundness*, but also onto, meaning *completeness*. Having short proofs means that the system is *polynomially bounded*: every tautology has a proof, i.e., preimage, of length polynomial in its length. Such proof systems exist if and only if  $\mathbf{NP} = \mathbf{coNP}$  [11].

As is well known, “one can think of length-of-proofs lower bounds as about problems of how to construct suitable models of particular bounded arithmetic” [23, p. 175]. A general method to construct such models is developed in [23] following Scott’s [35] forcing with random variables. An important instance of this method is the Boolean valued model  $K(F_{PV}^n)$ .<sup>1</sup> Its universe is given by the set of all polynomial time functions on binary strings of some fixed nonstandard length  $n \in M$ , where  $M$  is some fixed large nonstandard model of true arithmetic. The Boolean valuation considers two such functions equal if they differ only on an infinitesimal fraction of input strings. The model interprets the language having symbols for all polynomial time functions and relations, and it turns out that in  $K(F_{PV}^n)$  all true universal statements in this language are valid. In particular, this holds for Cook’s theory PV formalizing feasible reasoning [9]. This together with its appealing and familiar definition makes  $K(F_{PV}^n)$  an object of interest. We shall mention some related constructions (cf. Remark 3.3) once we gave the precise definition in § 3.2.

The objects of  $K(F_{PV}^n)$  “can be viewed from two different perspectives” [23, p. 160], namely, first as elements of the universe of  $K(F_{PV}^n)$  and second as functions defined on binary strings  $\{0, 1\}^n$ . For example, viewed as an element of  $K(F_{PV}^n)$  a propositional proof system is a tautology in the sense of  $K(F_{PV}^n)$ . Conversely, a tautology in the sense of  $K(F_{PV}^n)$  is a *pseudo proof system* (Definition 3.4). Viewed as a function on binary strings a pseudo proof system may be unsound. In fact, it is conceivable that *mad* pseudo proof systems exist (Definition 3.6). Viewed as elements of  $K(F_{PV}^n)$  these are tautologies in the sense of  $K(F_{PV}^n)$  but viewed as functions on binary strings they *never* output a tautology.

\* Corresponding author; e-mail: [jmaly@dbai.tuwien.ac.at](mailto:jmaly@dbai.tuwien.ac.at)

<sup>1</sup> A short introduction to the model-theoretic concepts is given in § 3, no knowledge of [23] is assumed.

In [23, § 24.4], Krajíček asks for transfer principles concerning pseudo proof systems. Loosely speaking, a *transfer principle* is a statement that allows to infer properties of standard objects from properties of nonstandard objects, and vice-versa. In this note we prove such a transfer principle that links the existence of mad pseudo proof systems to a hypothesis concerning the computational complexity of the satisfiability problem SAT of independent interest, explained next.

**Hard sequences.** For an algorithm solving a hard computational task there exist instances of the problem witnessing that the algorithm is not feasible. For example,  $\mathbf{P} \neq \mathbf{NP}$  if and only if every SAT-algorithm has a hard sequence:

**Definition 1.1** Let  $Q \subseteq \{0, 1\}^*$  and  $\mathbb{A}$  be a  $Q$ -algorithm, i.e., an algorithm deciding  $Q$ , and let  $p$  be a polynomial. A sequence  $(x_n)_{n \in \mathbb{N}}$  is  $p$ -hard for  $\mathbb{A}$  if for infinitely many  $n \in \mathbb{N}$

(H1)  $x_n \in Q$  and

(H2)  $t_{\mathbb{A}}(x_n) > p(|x_n|, n)$ .

Here,  $t_{\mathbb{A}}(x)$  denotes the running time of  $\mathbb{A}$  on input  $x$ . Being *hard for  $\mathbb{A}$*  means being  $p$ -hard for  $\mathbb{A}$  for all polynomials  $p$ .

It is a natural question to ask whether such a sequence could be computable in polynomial time. Here, we say that a sequence  $(x_n)_{n \in \mathbb{N}}$  of binary strings  $x_n \in \{0, 1\}^*$  is polynomial time computable if so is the function that computes  $x_n$  from  $1^n = 1 \cdots 1$  ( $n$  times).<sup>2</sup>

Hard sequences have been studied from at least two perspectives. The first is *speed-up*, going back at least to [36], and the second, more relevant to this paper, is to *witness failure* of feasible algorithms, studied not only for deterministic but also for randomized [3, 13, 37] and non-uniform algorithms [2, 27]. Cf. [6, 31] for some recent discussions.

**Hard Sequence Hypothesis 1.2** For every SAT-algorithm  $\mathbb{A}$  there exists a polynomial time computable sequence which is hard for  $\mathbb{A}$ .

We are not aware of a place where this hypothesis has been formulated explicitly, but it is certainly implicit in many papers. We are also not aware of any well-established computational hardness hypothesis that would imply this hypothesis.

A natural (cf., e.g., [13]) weakening of the Hard Sequence Hypothesis is to allow randomness in the construction of hard sequences. One then asks for polynomial time *samplable* (as opposed to computable) *probably hard* sequences (Definitions 2.2 & 2.1). We observe that SAT-algorithms do have such sequences under cryptographic assumptions (Proposition 2.4). Second, we observe that SAT-solvers do have polynomial time computable hard sequences under a well-established hypothesis (Proposition 2.7): SAT-algorithms which upon accepting a satisfiable input formula  $F$  also output a binary string of length  $\leq |F|$  that satisfies  $F$ . We say that  $x = x_1 \cdots x_n \in \{0, 1\}^n$  satisfies  $F$  if so does the truth assignment that maps the  $i$ -th variable of  $F$  to  $x_i$  if  $i \leq n$  and to 0 otherwise.

**Transfer principle.** Our transfer principle links the existence of mad pseudo proof systems with the existence of probably hard sequences that are samplable with a quite restrictive use of randomness that we call *invertibility* (cf. Definition 2.8). Intuitively, the sampler is required to witness its outputs by publishing the random seed used.

**Theorem 1.3** Let  $M$  be an  $\aleph_1$ -saturated elementary extension of  $\mathbb{N}$ . Then the following are equivalent:

- (a) There is a nonstandard  $n \in M$  such that  $K(F_{\text{PV}}^n)$  contains mad pseudo proof systems.
- (b) Every SAT-solver has an invertibly samplable probably hard sequence.
- (c) There is a polynomial time computable function  $f$  such that for all  $\varepsilon > 0$  and all polynomial time computable functions  $g$  there are infinitely many  $n \in \mathbb{N}$  such that
  - (i) for all  $x \in \{0, 1\}^n$ :  $f(x)$  is a falsifiable propositional formula;
  - (ii) for at most an  $\varepsilon$ -fraction of  $x \in \{0, 1\}^n$ :  $g(x)$  is a falsifying assignment of  $f(x)$ .

These statements are true if  $\mathbf{NE} \cap \mathbf{coNE} \not\subseteq \mathbf{E}$  and pseudo-random generators exist.

<sup>2</sup> Note that (H2) in the definition of a polynomial time computable hard sequence is equivalent to the statement that  $t_{\mathbb{A}}(x_n)$  is not  $n^{O(1)}$ , matching the definition in [8].

The property in (c) could be taken as a standard definition of a mad pseudo proof system, i.e., one not referring to nonstandard models. It is our notion of invertibility that makes the equivalence of (b) and (c) an easy consequence of Levin's optimal SAT-solver [26]. The equivalence to (a) is the remark this short paper wants to communicate, as a contribution to the model-theory of the structures  $K(F_{PV}^n)$ .

## 2 Hard sequences

Assuming  $\mathbf{NP} \not\subseteq \mathbf{P}$ , Gutfreund, Shaltiel, and Ta-Shma showed that for every fixed polynomial  $p$  every SAT-algorithm has polynomial time computable  $p$ -hard sequence [13]. A diagonalizing argument shows that one can compute hard sequences in slightly superpolynomial time (cf. [13, Theorem 1.6] for such a construction) but the construction of polynomial time computable hard sequences is open.

Does randomness help? As for a notion of feasibility for sequences of random strings<sup>3</sup> we borrow the following from average case complexity [4]:

**Definition 2.1** A sequence of random strings  $(X_n)_{n \in \mathbb{N}}$  is (*polynomial time*) *samplable* if there exists a polynomial time computable *sampler* for it, that is, a function  $D : \{0, 1\}^* \rightarrow \{0, 1\}^*$  such that  $D \circ U_n$  has the same distribution as  $X_n$  for all  $n \in \mathbb{N}$ . Here,  $U_n$  denotes a random variable uniformly distributed in  $\{0, 1\}^n$ .

The following definition is convenient. With suitable adjustments, it makes sense for *randomized* SAT-algorithms, and has been implicitly studied in [13, 37]. Here, we restrict attention to deterministic algorithms.

**Definition 2.2** Let  $\mathbb{A}$  be a  $Q$ -algorithm,  $p$  a polynomial and  $\delta, \varepsilon \geq 0$ . A sequence  $(X_n)_{n \in \mathbb{N}}$  of random strings is  $(\delta, \varepsilon)$ -*probably  $p$ -hard* for  $\mathbb{A}$  if for infinitely many  $n \in \mathbb{N}$

$$(P1) \Pr(X_n \in Q) \geq 1 - \delta \text{ and}$$

$$(P2) \Pr(t_{\mathbb{A}}(X_n) > p(|X_n|, n)) \geq 1 - \varepsilon.$$

The sequence is  $(\delta, \varepsilon)$ -*probably hard* for  $\mathbb{A}$  if for all polynomials  $p$  it is  $(\delta, \varepsilon)$ -probably  $p$ -hard for  $\mathbb{A}$ . And we call it *probably hard* for  $\mathbb{A}$  if for all  $\varepsilon > 0$  it is  $(0, \varepsilon)$ -probably hard for  $\mathbb{A}$ .

We observe that, using randomness, (superpolynomial) hardness is achievable under cryptographic assumptions:<sup>4</sup>

**Definition 2.3** A *cryptographic pseudo-random generator with stretch  $2n$*  is a polynomial time computable function  $G : \{0, 1\}^* \rightarrow \{0, 1\}^*$  such that  $|G(r)| = 2|r|$  for all  $r \in \{0, 1\}^*$  and for all positive polynomials  $p$  and all randomized polynomial time algorithms  $\mathbb{A}$  we have for all sufficiently large  $n$ :

$$|\Pr(\mathbb{A} \text{ accepts } G(U_n)) - \Pr(\mathbb{A} \text{ accepts } U_{2n})| \leq 1/p(n). \quad (1)$$

**Proposition 2.4** *Assume cryptographic pseudo-random generators with stretch  $2n$  exist. Then there is a samplable sequence which is probably hard for every SAT-algorithm.*

**Proof.** Let  $G$  be a generator as assumed to exist. Clearly, its image  $Q := \{G(r) \mid r \in \{0, 1\}^*\}$  is in  $\mathbf{NP}$ , so there is a polynomial time reduction  $f$  from  $Q$  to SAT. Define  $D(r) := f(G(r))$ , and note  $\Pr(D(U_n) \in \text{SAT}) = 1$  for all  $n \in \mathbb{N}$ . Assume for the sake of contradiction, that  $(D(U_n))_n$  is not probably hard for some SAT-algorithm  $\mathbb{B}$ . Then there are a polynomial  $p$  and  $\varepsilon > 0$  such that  $\Pr(t_{\mathbb{B}}(D(U_n)) \leq p(n)) \geq \varepsilon$  for infinitely many good  $n$ .

Let  $\mathbb{A}$  accept an input  $r$  if and only if  $\mathbb{B}$  accepts  $f(r)$  in at most  $p(|r|)$  steps. Then  $\Pr(\mathbb{A} \text{ accepts } G(U_n)) \geq \varepsilon$  for all good  $n$ . But the event that  $\mathbb{A}$  accepts  $U_{2n}$  implies the event that  $\mathbb{B}$  accepts  $f(U_{2n})$ , hence  $f(U_{2n}) \in \text{SAT}$ , hence  $U_{2n} \in Q$ . The latter event has probability  $\leq 2^n/2^{2n} = 2^{-n}$ . Thus, for all large enough good  $n$  the difference of the probabilities in (1) is at least  $\varepsilon - 2^{-n} \geq \varepsilon/2$ , a contradiction.  $\square$

While we are not aware of a reference for the above result, its proof is certainly folklore. Early references for similar constructions are [5, 10] where, instead of general SAT-algorithms, SAT-solvers are considered.

<sup>3</sup> A random string is a random variable with values in  $\{0, 1\}^*$ . Given any random variable we always use  $\Pr$  to denote the probability measure of the underlying probability space.

<sup>4</sup> Such assumptions are prohibitive in the context of [2, 3, 13] who are concerned with the problem to reduce average-case hardness hypotheses to worst-case hardness hypotheses.

Cryptographic assumptions are not required in this case, e.g., [3, 24] construct  $p$ -hard sequences assuming certain proof lower bounds and  $\mathbf{NP} \not\subseteq \mathbf{P}$  respectively. The sequences in [3, 5, 10] also produce formulas along with satisfying assignments, and it is clear that deterministic such sequences, namely *dreambreakers* [3], cannot be (superpolynomially) hard for Levin's optimal SAT-solver  $\mathbb{L}$ :

**Theorem 2.5** (Levin; [26]) *There exists a SAT-solver  $\mathbb{L}$  such that for every SAT-solver  $\mathbb{A}$  there exists a polynomial  $p_{\mathbb{A}}$  such that  $t_{\mathbb{L}}(F) \leq p_{\mathbb{A}}(t_{\mathbb{A}}(F) + |F|)$  for every  $F \in \text{SAT}$ .*

We shall use the following easy consequence mainly with  $\delta = 0$  by referring to “the optimality of  $\mathbb{L}$ ”. Note the lemma applies to (deterministic) hard sequences because these are  $(0,0)$ -probably hard sequences.

**Lemma 2.6** *Let  $(X_n)_{n \in \mathbb{N}}$  be a sequence of random strings and  $\varepsilon \geq \delta \geq 0$ . If  $(X_n)_{n \in \mathbb{N}}$  is  $(\delta, \varepsilon - \delta)$ -probably hard for  $\mathbb{L}$ , then it is  $(\delta, \varepsilon)$ -probably hard for every SAT-solver.*

**Proof.** Let  $\mathbb{A}$  be a SAT-solver. Assume for the sake of a contradiction that  $(X_n)_{n \in \mathbb{N}}$  is  $(\delta, \varepsilon - \delta)$ -probably hard for  $\mathbb{L}$  but not for  $\mathbb{A}$ . Then, let  $p$  be a polynomial such that for almost all  $n \in \mathbb{N}$ :

$$\Pr(X_n \in \text{SAT}) \geq 1 - \delta \text{ implies } \Pr(t_{\mathbb{A}}(X_n) \leq p(|X_n|, n)) \geq \varepsilon. \quad (2)$$

Choose a nondecreasing polynomial  $p_{\mathbb{A}}$  for  $\mathbb{A}$  according to Theorem 2.5. Then (2) implies

$$\Pr(X_n \in \text{SAT}) \geq 1 - \delta \text{ implies } \Pr(t_{\mathbb{L}}(X_n) \leq p_{\mathbb{A}}(p(|X_n|, n)) \text{ or } X_n \notin \text{SAT}) \geq \varepsilon,$$

and thus

$$\Pr(X_n \in \text{SAT}) \geq 1 - \delta \text{ implies } \Pr(t_{\mathbb{L}}(X_n) \leq p_{\mathbb{A}}(p(|X_n|, n))) \geq \varepsilon - \delta.$$

Hence,  $(X_n)_{n \in \mathbb{N}}$  is not  $(\delta, \varepsilon - \delta)$ -probably hard for  $\mathbb{L}$ . Contradiction.  $\square$

We point out that one can construct superpolynomially hard sequences for SAT-solvers under standard worst-case assumptions. The proof is essentially known. Recall,  $\mathbf{E}$  and  $\mathbf{NE}$  denote deterministic and nondeterministic simply exponential time  $2^{O(n)}$ , respectively.

**Proposition 2.7** *The following statements are equivalent, and implied by  $\mathbf{NE} \cap \mathbf{coNE} \not\subseteq \mathbf{E}$ :*

- (a) *There exists a polynomial time computable sequence which is hard for  $\mathbb{L}$ .*
- (b) *There exists a polynomial time computable sequence which is hard for all SAT-solvers.*
- (c) *For every SAT-solver  $\mathbb{A}$  there exists a polynomial time computable sequence which is hard for  $\mathbb{A}$ .*
- (d) *For every SAT-solver  $\mathbb{A}$  there exists an injective polynomial time computable sequence  $(F_n)_{n \in \mathbb{N}}$  which is hard for  $\mathbb{A}$  and such that  $|F_n| \geq n$  for all  $n \in \mathbb{N}$ .*

**Proof.** (a) $\Rightarrow$ (b) follows from the optimality of  $\mathbb{L}$  (Lemma 2.6). (b) $\Rightarrow$ (c) and (d) $\Rightarrow$ (a) are trivial. To prove (c) $\Rightarrow$ (d) we proceed as in [8, Proposition 3.2] using a padding function: a polynomial time computable function  $\text{pad}$  that maps a formula  $F$  and a string  $y \in \{0, 1\}^*$  to a formula  $\text{pad}(F, y)$  of length at least  $|F| + |y|$  that has the same satisfying assignments as  $F$ , and such that there are two polynomial time functions mapping any input of the form  $\text{pad}(F, y)$  to  $F$  and  $y$ , respectively.

Let  $\mathbb{A}$  be a SAT-solver and assume (c). Define an algorithm  $\mathbb{B}$  as follows: given a formula  $F$ , for  $t = 0, 1, \dots$  compute  $t$  steps of  $\mathbb{A}$  on each of  $\text{pad}(F, 1^0), \dots, \text{pad}(F, 1^t)$ ; as soon as one of these computations halts, return the answer obtained.

Clearly,  $\mathbb{B}$  is a SAT-solver and there is a polynomial  $p$  such that for every  $t \in \mathbb{N}$  and every formula  $F$  we have  $t_{\mathbb{B}}(F) \leq p(t + t_{\mathbb{A}}(\text{pad}(F, 1^t)))$ . By (c) there is a polynomial time computable sequence  $(F_n)_n$  hard for  $\mathbb{B}$ . Then  $(\text{pad}(F_n, 1^n))_n$  is polynomial time computable and hard for  $\mathbb{A}$ . This sequence is injective and satisfies  $|\text{pad}(F_n, 1^n)| \geq n$  for all  $n \in \mathbb{N}$ .

We have proved that (a)–(d) are equivalent. We now derive (b) assuming there exists a problem  $Q \in \mathbf{NE} \cap \mathbf{coNE} \setminus \mathbf{E}$ . For a binary string  $x$  let  $\text{num}(x)$  be the natural number with binary expansion  $1x$ . Then

$$Q' := \{1^{\text{num}(x)} \mid x \in Q\} \in \mathbf{NP} \cap \mathbf{coNP} \setminus \mathbf{P}.$$

We now proceed as in [7, Proposition 4.5]. By the  $\mathbf{NP}$ -completeness of SAT, there are polynomial time reductions  $r_1$  and  $r_0$  from  $Q'$  and  $\{0, 1\}^* \setminus Q'$  to SAT. We can assume that  $r_1(1^n)$  and  $r_0(1^n)$  are propositional formulas. Then

$r_1(1^n) \vee r_0(1^n) \in \text{SAT}$ , and a satisfying assignment satisfies exactly one of  $r_1(1^n)$  and  $r_0(1^n)$ , namely  $r_1(1^n)$  if  $1^n \in Q'$ , and  $r_0(1^n)$  if  $1^n \notin Q'$ . Since  $Q' \notin \mathbf{P}$ , there is no SAT-solver  $\mathbb{A}$  such that  $t_{\mathbb{A}}(r_1(1^n) \vee r_0(1^n)) \leq n^{O(1)}$ .  $\square$

We now consider sequences sampled with some restricted use of randomness, as announced § 1.<sup>5</sup>

**Definition 2.8** A sequence of random strings  $(X_n)_{n \in \mathbb{N}}$  is *invertibly samplable* if it has a polynomial time sampler  $D$  which is *invertible*, i.e.,  $D$  is injective and the partial function  $D^{-1}$  is computable in polynomial time.

The sampler defined in the proof of Proposition 2.4 is not invertible. For invertible samplers, hardness has the following handy reformulation.

**Lemma 2.9** *Let  $D$  be an invertible polynomial time sampler for  $(X_n)_{n \in \mathbb{N}}$ . Then the following are equivalent:*

- (a)  $(X_n)_{n \in \mathbb{N}}$  is probably hard for  $\mathbb{L}$ .
- (b) For every polynomial time function  $g$  and for all  $\varepsilon > 0$  there are infinitely many  $n$  such that  $\Pr(X_n \in \text{SAT}) = 1$  and

$$\Pr\left(|g(U_n)| \leq |D(U_n)| \text{ and } g(U_n) \text{ satisfies } D(U_n)\right) \leq \varepsilon. \quad (3)$$

**Proof.** “(a) $\Rightarrow$ (b)”. Assume (b) fails and choose  $g$  and  $\varepsilon$  witnessing this. Define the following algorithm  $\mathbb{A}$ : given as input a formula  $F$ , compute the string  $y := g(D^{-1}(F))$  and check whether it has length  $\leq |F|$  and satisfies  $F$ ; if so, then accept with output  $y$ , else reject.

Further, define the algorithm  $\mathbb{B}$  to run  $\mathbb{A}$  in parallel with an arbitrary SAT-solver. If one of the two procedures halts accepting, then  $\mathbb{B}$  accepts with the corresponding output. If both procedures reject, so does  $\mathbb{B}$ .

Since  $\mathbb{A}$  is polynomial time bounded, there is a polynomial  $p$  such that  $t_{\mathbb{A}}(F) \leq p(|F|)$  for every formula  $F$  accepted by  $\mathbb{A}$ . Since  $\mathbb{B}$  is a SAT-solver, there exists a polynomial  $p_{\mathbb{B}}$  such that then  $t_{\mathbb{L}}(F) \leq p_{\mathbb{B}}(p(|F|))$  (Theorem 2.5). Thus

$$\begin{aligned} \Pr\left(t_{\mathbb{L}}(X_n) \leq p_{\mathbb{B}}(p(|X_n|))\right) &\geq \Pr\left(\mathbb{A} \text{ accepts } X_n\right) \\ &= \Pr\left(|g(D^{-1}(X_n))| \leq |X_n| \text{ and } g(D^{-1}(X_n)) \text{ satisfies } X_n\right). \end{aligned}$$

Note this last probability equals the probability in (3). By our assumption that (a) fails this probability is  $> \varepsilon$  or  $\Pr(X_n \in \text{SAT}) < 1$  for almost all  $n$ . Hence,  $(X_n)_{n \in \mathbb{N}}$  is not  $(0, \varepsilon)$ -probably  $(p_{\mathbb{B}} \circ p)$ -hard for  $\mathbb{L}$ .

“(b) $\Rightarrow$ (a)”. If (a) fails, there is a polynomial  $p$  and an  $\varepsilon > 0$  such that for almost all  $n$ ,  $\Pr(X_n \in \text{SAT}) < 1$  or  $\Pr(t_{\mathbb{L}}(X_n) \leq p(|X_n|, n)) > \varepsilon$ .

Define a polynomial time function  $g$  as follows. On input  $r$  run  $\mathbb{L}$  on  $D(r)$  for at most  $p(|D(r)|, |r|)$  steps. If this computation does not halt accepting, then return the empty string; else return  $\mathbb{L}$ 's output. Then  $|g(r)| \leq |D(r)|$  for all  $r$ , and the probability in (3) equals the probability of the event that  $g(U_n)$  satisfies  $D(U_n)$ . For  $n$  with  $\Pr(X_n \in \text{SAT}) = 1$ , this event is implied by the event that  $t_{\mathbb{L}}(X_n) \leq p(|X_n|, n)$ , so has probability  $> \varepsilon$ . Thus statement (b) fails.  $\square$

An easy corollary is the equivalence of statement (b) and (c) of our main theorem.

**Proof of Theorem 1.3 (b) $\Leftrightarrow$ (c).** It follows from the optimality of  $\mathbb{L}$  that there is  $D$  as in Lemma 2.9 (a) if and only if Theorem 1.3 (b) holds.

We show that there is  $D$  as in Lemma 2.9 (b) if and only if Theorem 1.3 (c) holds. Given a sampler  $D$  as in Lemma 2.9 (b), define  $f(r) := \text{neg}(D(r))$ , where  $\text{neg}$  maps (a binary string coding) a formula to a its negation. Conversely, given a function  $f$  as in Theorem 1.3 (c), define the invertible sampler  $D(r) := \text{pad}(\text{neg}(f(r)), r)$ , where  $\text{pad}$  is the padding function from the proof of Proposition 2.7.

For every  $n \in \mathbb{N}$  and every polynomial time computable function  $g$  we have:  $D(r)$  is a satisfiable formula for all  $r \in \{0, 1\}^n$  if and only if  $f(r)$  is a falsifiable formula for all  $r \in \{0, 1\}^n$ ; further, the probability that  $g(U_n)$  satisfies  $D(U_n)$  is  $\leq \varepsilon$  if and only if  $g(r)$  is a falsifies  $f(r)$  for only a  $\varepsilon$ -fraction of  $r \in \{0, 1\}^n$ .  $\square$

<sup>5</sup> Our notion of invertibility is more restrictive than the one considered in [38].

Hard sequences can be transformed into invertibly samplable probably hard sequences using pseudo-random generators (of the Nisan-Wigderson type). This is a standard application of the “general framework for derandomization” of [21]. For definiteness we use the parameter setting from the standard textbook [1].

**Definition 2.10** Let  $S : \mathbb{N} \rightarrow \mathbb{N}$ . A function  $G : \{0, 1\}^* \rightarrow \{0, 1\}^*$  is an  $S(\ell)$ -pseudo-random generator if  $G(r)$  is computable in time  $2^{O(|r|)}$ , has length  $S(|r|)$  and for all  $\ell \in \mathbb{N}$  and all Boolean circuits  $C$  with at most  $S(\ell)^3$  gates and at most  $S(\ell)$  inputs

$$|\Pr(C(G(U_\ell)) = 1) - \Pr(C(U_{S(\ell)}) = 1)| < 0.1. \tag{4}$$

We say pseudo-random generators exist if there is  $\delta > 0$  such that  $2^{\lfloor \delta \ell \rfloor}$ -pseudo-random generators exist.

**Proposition 2.11** Assume pseudo-random generators exist. If there exists a polynomial time computable hard sequence for  $\mathbb{L}$ , then there exists an invertibly samplable probably hard sequence for  $\mathbb{L}$ .

*Proof.* Let  $(F_n)_n$  be polynomial time computable and hard for  $\mathbb{L}$ . By Proposition 2.7 we can assume that the sequence is injective and  $|F_n| \geq n$  for all  $n$ . Using the padding function  $\text{pad}$  from the proof of this proposition, define a sampler  $D(r) := \text{pad}(F_{|r|}, r)$ .

Clearly,  $D$  is polynomial time computable and invertible. Assume for the sake of contradiction that  $(D(U_n))_n$  is not probably hard for  $\mathbb{L}$ . Applying Lemma 2.9 we get a polynomial time function  $g$  and  $\varepsilon > 0$  and  $n_0 \in \mathbb{N}$  such that for all  $n > n_0$

$$\Pr(D(U_n) \in \text{SAT}) = 1 \text{ implies } \Pr(|g(U_n)| \leq |D(U_n)| \text{ and } g(U_n) \text{ satisfies } D(U_n)) > \varepsilon. \tag{5}$$

Note  $\Pr(D(U_n) \in \text{SAT})$  is 1 or 0 depending on whether  $F_n \in \text{SAT}$  or not. Further note that a string satisfies  $D(U_n)$  if and only if it satisfies  $F_n$ . Hence (5) implies

$$F_n \in \text{SAT} \text{ implies } \Pr(g(U_n) \text{ satisfies } F_n) > \varepsilon. \tag{6}$$

Call  $n \in \mathbb{N}$  good if  $n > n_0$  and  $F_n \in \text{SAT}$ . We claim there is a SAT-solver  $\mathbb{A}$  such that  $t_{\mathbb{A}}(F_n) \leq n^{O(1)}$  for all good  $n$ . This implies that  $(F_n)_n$  is not hard for  $\mathbb{A}$  and thus also not for  $\mathbb{L}$  (Lemma 2.6), a contradiction.

Let  $c \in \mathbb{N}$  be such that  $(1 - \varepsilon)^c \leq 0.9$ . Let  $C_n$  be a size  $n^{O(1)}$  circuit with  $c \cdot n$  inputs that accepts  $r_1 \cdots r_c$  with  $r_i \in \{0, 1\}^n$  if and only if at least one of  $g(r_1), \dots, g(r_c)$  satisfies  $F_n$ . If  $n$  is good, then  $\Pr(C_n(U_{cn}) = 1) > 0.1$  by choice of  $c$  and (6). For every  $m \geq cn$  we can view  $C_n$  as a circuit  $C'_m$  on  $m$  inputs. Further, there is a polynomial time function  $g$  such that  $g(n, r)$  satisfies  $F_n$  whenever  $r \in \{0, 1\}^m$  is such that  $C'_m(r) = 1$ .

If we set  $m_n := 2^{\lfloor \delta \ell_n \rfloor}$  where  $\ell_n := \lfloor d \log n \rfloor$  for a sufficiently large constant  $d \in \mathbb{N}$ , then  $m_n \geq cn$  and  $C'_{m_n}$  has size  $\leq m_n^3$ . Here,  $\delta > 0$  witnesses that there exists a pseudo-random generator  $G$ . For all good  $n$  we have  $\Pr(C'_{m_n}(U_{m_n}) = 1) > 0.1$ , so  $\Pr(C'_{m_n}(G(U_{\ell_n})) = 1) \neq 0$  by (4). Hence, for good  $n$ ,  $g(n, G(r))$  satisfies  $F_n$  for at least one  $r \in \{0, 1\}^{\ell_n}$ .

Define the SAT-solver  $\mathbb{A}$  as follows. Given a formula  $F$  it runs some arbitrary SAT-solver and in parallel does the following: compute  $F_0, \dots, F_{|F|}$ ; unless there is  $n_0 < n \leq |F|$  such that  $F_n = F$ , reject; otherwise compute the strings  $g(n, G(r))$  for all  $\leq n^d$  many  $r \in \{0, 1\}^{\ell_n}$ ; if one of them satisfies  $F = F_n$ , then output it and accept; else reject.

It is easy to see that  $\mathbb{A}$  is polynomially time bounded on  $F_n$  for good  $n$ , as desired. □

### 3 Mad pseudo proof systems

#### 3.1 Preliminaries: language $L_{PV}$ and model $M$

So far we considered polynomial time on the set of binary strings  $\{0, 1\}^*$ . To view polynomial time on  $\mathbb{N}$ , we view every  $n \in \mathbb{N}$  as a binary string, say, by taking the binary expansion of  $n$  and deleting the most significant bit. Then  $\{0, 1\}^n$  corresponds to the numbers between  $2^n$  and  $2^{n+1} - 1$ , and we continue to write  $\{0, 1\}^n$  for this interval.

We consider every  $r$ -ary polynomial time computable function  $f : \mathbb{N}^r \rightarrow \mathbb{N}$  as an  $r$ -ary function symbol and every  $r$ -ary polynomial time decidable relation  $R \subseteq \mathbb{N}^r$  as an  $r$ -ary relation symbol. Constants are nullary function symbols. Let  $L_{PV}$  denote the resulting first-order language. The standard  $L_{PV}$ -structure has as universe  $\mathbb{N}$  and interprets all function and relation symbols from  $L_{PV}$  by themselves. We denote this structure also by  $\mathbb{N}$  and in

general do not distinguish structures from their universes notationally. The theory  $\text{Th}_V(L_{PV})$  is the set of universal sentences true in  $\mathbb{N}$ .

To fix some notation we list some symbols of the language  $L_{PV}$ . It contains a unary function  $|n|$  denoting the length of  $n$  as a binary string; the  $\{0, 1\}$ -valued binary function  $\text{bit}(i, n)$  gives the  $i$ -th bit of this string, and  $0$  if  $i > |n|$ . We say  $n$  codes the set  $\sqcup n \sqcup := \{i \in \mathbb{N} \mid \text{bit}(i, n) = 1\}$ . For a set  $A \subseteq \mathbb{N}$  coded in  $\mathbb{N}$  let  $\ulcorner A \urcorner$  denote its code. A finite function  $\alpha$  is coded by  $m$  if  $m$  codes the set of  $\langle i, \alpha(i) \rangle$  for  $i$  in the domain of  $\alpha$ ; here,  $\langle i, j \rangle$  is a bijection from  $\mathbb{N}^2$  onto  $\mathbb{N}$ . For readability we write  $i \in A$  for  $\text{bit}(i, \ulcorner A \urcorner) = 1$ , and  $\alpha(i)$  for a suitable  $L_{PV}$ -term applied to  $i$  and the code of  $\alpha$ .

Positive rationals are coded by pairs  $\langle n, m \rangle$  written  $n/m$  with  $m \neq 0$  and ambiguously we use the symbol  $\leq$  also with its meaning in the rationals. There is a unary function  $\text{card}(n)$  in  $L_{PV}$  giving the cardinality of  $\sqcup n \sqcup$ . Further,  $L_{PV}$  contains a unary function mapping  $\ulcorner \emptyset \urcorner$  to  $0$ , and  $\ulcorner A \urcorner$  to the rational  $\text{card}(\ulcorner A \urcorner)/2^n = \text{Pr}(U_n \in A)$  for every nonempty  $A \subseteq \{0, 1\}^n$ . We also write  $\text{Pr}$  for this function.

We fix an  $\aleph_1$ -saturated elementary extension  $M$  of  $\mathbb{N}$ . This means  $M$  is an extension of  $\mathbb{N}$  with the property that every countable family of definable subsets of  $M$  with the finite intersection property has non-empty intersection. By *definable* we mean definable by formulas with parameters from  $M$ . Elements of  $M \setminus \mathbb{N}$  are *nonstandard*.

We speak of sets and functions coded in  $M$  in the same sense as explained above, in particular, we have the notations  $\sqcup a \sqcup$  and  $\ulcorner A \urcorner$  for elements  $a$  of  $M$  and coded subsets  $A$  of  $M$ . The interpretation of a symbol  $\sigma \in L_{PV}$  in  $M$  is denoted  $\sigma^M$  but we shall often omit the superscript  $M$ . Pairs  $\langle a, b \rangle = \langle a, b \rangle^M$  written  $a/b$  with  $b \neq 0$  are *M-rationals*. E.g., the values of  $\text{Pr}^M$  are *M-rationals*. Note that every (code of a) rational is an *M-rational*.

We use the following notions from nonstandard analysis (cf., e.g., [19]). The *standard part* of an *M-rational*  $a/b$  is the real

$$(a/b)^* := \inf\{q \in \mathbb{Q} \mid a/b \leq q\},$$

provided the set on the r.h.s. is non-empty; it is undefined otherwise. An *M-rational* with standard part  $0$  is *infinitesimal*.

### 3.2 Krajíček's model $K(F_{PV}^n)$

The model  $K(F_{PV}^n)$  is Boolean valued with values in the Boolean algebra  $\mathcal{B}_n$ , defined below.

The function  $n \mapsto \ulcorner \{0, 1\}^n \urcorner$  is definable in the standard model  $\mathbb{N}$ . Since  $M$  is an elementary extension of  $\mathbb{N}$ , this function extends to a function on  $M$ . Evaluating it on  $n \in M$  gives the code of a subset of  $M$  that we denote by  $\{0, 1\}^n$ . Let  $\mathcal{A}_n$  be the set of subsets of  $\{0, 1\}^n$  that are coded in  $M$ . Then  $\mathcal{A}_n$  is a Boolean algebra and  $\{\ulcorner A \urcorner \mid A \in \mathcal{A}_n\}$  is coded in  $M$ . Note that for every  $L_{PV}$ -formula  $\varphi(x)$  (even with parameters from  $M$ ) we have  $\{\omega \in \{0, 1\}^n \mid M \models \varphi(\omega)\} \in \mathcal{A}_n$ .

Let  $n \in M$  be nonstandard. The set  $\text{Inf}_n := \{A \in \mathcal{A}_n \mid \text{Pr}^M(\ulcorner A \urcorner) \text{ is infinitesimal}\}$  is an ideal in  $\mathcal{A}_n$  (and not coded in  $M$ ). Call  $A, A'$  equivalent if their symmetric difference is in  $\text{Inf}_n$ . The equivalence class of  $A \in \mathcal{A}_n$  is denoted  $A/\text{Inf}_n$ . These classes form the Boolean algebra  $\mathcal{B}_n$ , defined as the factor  $\mathcal{B}_n := \mathcal{A}_n/\text{Inf}_n$ .

Using the assumption that  $M$  is  $\aleph_1$ -saturated one can show [23, Lemma 1.2.1]:

**Lemma 3.1** *For every nonstandard  $n \in M$ , the Boolean algebra  $\mathcal{B}_n$  is complete.*

We now describe the model  $K(F_{PV}^n)$ . Its universe is  $F_{PV}^n$ , the set of all restrictions  $f^M \upharpoonright \{0, 1\}^n$  of  $f^M$  to  $\{0, 1\}^n$  where  $f \in L_{PV}$  is a unary function symbol (and  $f^M$  its interpretation in  $M$ ). We use  $\alpha, \beta, \dots$  to range over  $F_{PV}^n$ . Observe that every  $\alpha \in F_{PV}^n$  is coded in  $M$  (but not the set  $F_{PV}^n$ , viewed as a set of codes). Further observe that for every  $r$ -ary symbol  $f \in L_{PV}$  and every  $r$ -tuple  $(\alpha_1, \dots, \alpha_r)$  the function  $\omega \mapsto f^M(\alpha_1(\omega), \dots, \alpha_r(\omega))$  defined on  $\omega \in \{0, 1\}^n$  is in  $F_{PV}^n$ . We interpret the function symbols of  $L_{PV}$  in this way over  $F_{PV}^n$ . Then every closed  $L_{PV}$ -term  $t$  with parameters from  $F_{PV}^n$  denotes an element  $t^{K(F_{PV}^n)}$  of  $F_{PV}^n$ . The Boolean valuation maps every  $L_{PV}$ -sentence  $\varphi$  with parameters from  $F_{PV}^n$  to a Boolean value  $\llbracket \varphi \rrbracket \in \mathcal{B}_n$ . For atomic  $\varphi$  this Boolean value is defined setting:

$$\begin{aligned} \llbracket R(t_1, \dots, t_r) \rrbracket &:= \left\{ \omega \in \{0, 1\}^n \mid \left( t_1^{K(F_{PV}^n)}(\omega), \dots, t_r^{K(F_{PV}^n)}(\omega) \right) \in R^M \right\} / \text{Inf}_n, \\ \llbracket t = s \rrbracket &:= \left\{ \omega \in \{0, 1\}^n \mid t^{K(F_{PV}^n)}(\omega) = s^{K(F_{PV}^n)}(\omega) \right\} / \text{Inf}_n, \end{aligned}$$

where  $t, s, t_1, \dots, t_r$  are closed  $L_{PV}$ -terms with parameters from  $F_{PV}^n$  and  $R \in L_{PV}$  is an  $r$ -ary relation symbol. For arbitrary sentences with parameters in  $F_{PV}^n$  the Boolean value is then determined via the usual recurrence:  $\llbracket \neg\varphi \rrbracket := \sim\llbracket \varphi \rrbracket$ ,  $\llbracket (\varphi \vee \psi) \rrbracket := \llbracket \varphi \rrbracket \cup \llbracket \psi \rrbracket$ ,  $\llbracket \exists x\varphi(x) \rrbracket := \sup_\alpha \llbracket \varphi(\alpha) \rrbracket$  where  $\sim, \cup, \sup$  denote the obvious operations of  $\mathcal{B}_n$  as a complete Boolean algebra. The minimal and maximal elements of  $\mathcal{B}_n$  are respectively  $0_{\mathcal{B}_n} := \emptyset/\text{Inf}_n$  and  $1_{\mathcal{B}_n} := \{0, 1\}^n/\text{Inf}_n$ .

A sentence  $\varphi$  is *valid in*  $K(F_{PV}^n)$  if  $\llbracket \varphi \rrbracket = 1_{\mathcal{B}_n}$ . One straightforwardly verifies [23, Lemma 1.4.2]:

**Lemma 3.2** *Let  $n \in M$  be nonstandard. If  $\varphi(x, y, \dots)$  is a quantifier-free  $L_{PV}$ -formula and  $\alpha, \beta, \dots \in F_{PV}^n$ , then*

$$\llbracket \varphi(\alpha, \beta, \dots) \rrbracket = \{\omega \in \{0, 1\}^n \mid M \models \varphi(\alpha(\omega), \beta(\omega), \dots)\} / \text{Inf}_n.$$

*In particular, every sentence in  $\text{Th}_\forall(L_{PV})$  is valid in  $K(F_{PV}^n)$ .*

We close this subsection with some historical notes meant to back up our claim from the Introduction that the definition of  $K(F_{PV}^n)$  follows natural and familiar lines.

**Remark 3.3** (Historical notes) Boolean valued models date back to the work of Rasiowa and Sikorski [34], and became popular when it was realized that Cohen's method of forcing can be viewed as a method to construct Boolean valued models of set theory. We refer to [17, Chapter 14] and the references therein. In [35] Scott explained this view by constructing a model based on random variables of a higher-order theory of the reals as an ordered field. Such Boolean powers are studied in more generality in [29, 32].

The book [23] develops Scott's [35] forcing with random variables as a method to build models  $K(F)$  (and two-sorted extensions thereof) of bounded arithmetics. Instead of  $F_{PV}^n$  these use suitable families  $F \subseteq M^\Omega$  for  $\Omega$  coded in  $M$ , together with an analogously defined complete Boolean algebra  $\mathcal{B}$ . The crucial move being to restrict the construction to families  $F$  of random variables samplable with limited computational complexity. Technically, fullness ([17, p. 208], [29, Theorem 1.4]) of the model is lost and much of the theory develops around finding conditions ensuring *partial* fullness for certain classes of formulas.

The models  $K(F)$  can be seen as *partial* randomizations of  $M$  in the sense of Keisler [20]: the triple  $(F, \text{Pr}^M, \mathcal{B})$  satisfies only a fragment of Keisler's randomization theory. In particular,  $K(F)$  satisfies Keisler's "Fullness Axiom" [20, p. 128] only for very special  $F$  (cf. [23, Theorem 3.5.2]), and  $K(F_{PV}^n)$  does not.

As remarked in [23, footnote 2, p. 3] one can collapse  $K(F_{PV}^n)$  to a usual two-valued model by factoring  $\mathcal{B}_n$  with a suitable ultrafilter (cf. [34, Lemma 9.1]). The result is a restricted ultrapower of  $M$ . These have been studied for fragments of arithmetic [12, 16, 22, 25, 30, 33] ever since Skolem's definable ultrapower (cf. [14, IV.1.(b)]).

### 3.3 Pseudo proof systems

Let  $\text{Fml}$  be the set of naturals which (viewed as binary strings) code propositional formulas, and  $\text{Sat}$  contain the pairs  $(\ell, m)$  such that  $m \in \text{Fml}$  and  $\ell$  (as a binary string) satisfies the formula coded by  $m$ . Then  $\text{Fml}$  and  $\text{Sat}$  are relation symbols in  $L_{PV}$ . The formula

$$\text{Taut}(x) := \forall y (|y| \leq |x| \rightarrow \text{Sat}(y, x))$$

defines  $\text{TAUT}$ , viewed as a set of naturals. It follows from Lemma 3.2 that, if  $f \in L_{PV}$  is a proof system (i.e.,  $\forall x \text{Taut}(f(x)) \in \text{Th}_\forall(L_{PV})$ ), then  $f^M \upharpoonright \{0, 1\}^n \in F_{PV}^n$  is a pseudo proof system as defined in [23, p. 162]:

**Definition 3.4** Let  $n \in M$  be nonstandard. An element  $\alpha \in F_{PV}^n$  is a *pseudo proof system in*  $K(F_{PV}^n)$  if  $\text{Taut}(\alpha)$  is valid in  $K(F_{PV}^n)$ .

Hirsch, Itsykson, Monakhov, and Smal study *heuristic proof systems* in [15]. These are randomized proof systems that are allowed to prove non-tautologies (with constant probability) but only *few* of them with respect to some distribution. Pseudo proof systems are conceptually different. First, they are not randomized. More importantly, the point of a pseudo proof system is that erroneous outputs (non-tautologies) are hard to detect as such, and not that there are few of them. In fact, as we shall see in the next section, it is conceivable that there are *mad* pseudo proof systems, pseudo proof systems all of whose outputs are erroneous. The notion of a pseudo proof system is more akin to Kabanets's  $\text{pseudo}_P$ -classes [18].



Let  $\text{neg} \in L_{PV}$  map every  $n \in \text{Fml}$  to its negation (to the number coding the negation of the formula coded by  $n$ ).

**Lemma 3.5** *Let  $n \in M$  be nonstandard and  $f \in L_{PV}$ . The following are equivalent.*

- (a)  $f^M \upharpoonright \{0, 1\}^n$  is a pseudo proof system in  $K(F_{PV}^n)$ .  
 (b) For all  $g \in L_{PV}$ :

$$\{\omega \in \{0, 1\}^n \mid M \models |g(\omega)| \leq |f(\omega)| \wedge \neg \text{Sat}(g(\omega), f(\omega))\} \in \text{Inf}_n.$$

If  $M \models \forall x \in \{0, 1\}^n \text{Fml}(f(x))$ , then these statements are equivalent to

- (c) For all  $g \in L_{PV}$ :

$$\{\omega \in \{0, 1\}^n \mid M \models \text{Sat}(g(\omega), \text{neg}(f(\omega)))\} \in \text{Inf}_n.$$

**Proof.** Write  $\alpha := f^M \upharpoonright \{0, 1\}^n$ . Let  $\beta$  range over  $F_{PV}^n$ . Statement (a) means

$$\begin{aligned} 0_{B_n} &= \sup_{\beta} \sim \llbracket |\beta| \leq |\alpha| \rightarrow \text{Sat}(\beta, \alpha) \rrbracket \\ &= \sup_{\beta} \sim \{\omega \in \{0, 1\}^n \mid M \models |\beta(\omega)| \leq |\alpha(\omega)| \rightarrow \text{Sat}(\beta(\omega), \alpha(\omega))\} / \text{Inf}_n, \end{aligned}$$

using Lemma 3.2. Equivalently, for all  $\beta$ :

$$\{\omega \in \{0, 1\}^n \mid M \models |\beta(\omega)| \leq |f(\omega)| \wedge \neg \text{Sat}(\beta(\omega), f(\omega))\} \in \text{Inf}_n.$$

This is equivalent to (b) because the  $\beta \in F_{PV}^n$  are precisely the functions of the form  $g \upharpoonright \{0, 1\}^n$  for  $g \in L_{PV}$ .

Suppose  $M \models \forall x \in \{0, 1\}^n \text{Fml}(f(x))$ . Then for all  $g \in L_{PV}$ ,  $\text{Sat}(g(x), \text{neg}(f(x)))$  is equivalent to  $\neg \text{Sat}(g(\omega), f(\omega))$  in  $M$ , so (c) implies (b). Conversely, given  $g \in L_{PV}$  there is  $g' \in L_{PV}$  such that the set in (c) for  $g$  equals the set in (b) for  $g'$ : if  $|g(\omega)| > |f(\omega)|$ , then  $g'(\omega)$  deletes the last  $|g(\omega)| - |f(\omega)|$  many bits; this truncation does not change how the variables of the formula  $f(\omega)$  are evaluated.  $\square$

We are interested in pseudo proof systems that never output a tautology.

**Definition 3.6** Let  $n \in M$  be nonstandard. A pseudo proof system  $f^M \upharpoonright \{0, 1\}^n$  in  $K(F_{PV}^n)$  is *mad* if

$$M \models \forall x \in \{0, 1\}^n (\text{Fml}(f(x)) \wedge \neg \text{Taut}(f(x))). \quad (7)$$

### 3.4 Proof of Theorem 1.3

We already showed that (b) is equivalent to (c) and now show that (a) is equivalent to (c).

“(a) $\Rightarrow$ (c)”. Suppose there are  $n \in M \setminus \mathbb{N}$  and  $f \in L_{PV}$  such that  $f^M \upharpoonright \{0, 1\}^n$  is a mad pseudo proof system. We claim that  $f^{\mathbb{N}}$  satisfies (i) and (ii) of (c). By (7) we have  $M \models \psi(n)$  where  $\psi(x)$  is the  $L_{PV}$ -formula

$$\forall y \in \{0, 1\}^x \exists z \text{Sat}(z, \text{neg}(f(y))). \quad (8)$$

Now, let  $A_{g,f}(x)$  be the function, definable in  $\mathbb{N}$ , mapping  $m$  to the code of

$$\{\omega \in \{0, 1\}^m \mid \mathbb{N} \models |g(\omega)| \leq |\text{neg}(f(\omega))| \wedge \text{Sat}(g(\omega), \text{neg}(f(\omega)))\}.$$

Since  $M \models \forall x \in \{0, 1\}^n \text{Fml}(f(x))$  by (7), we get  $\perp A_{g,f}^M(n) \perp \in \text{Inf}_n$  by Lemma 3.5 (c). This implies for all  $g \in L_{PV}$  and all standard  $\ell > 0$  that  $M$  and hence  $\mathbb{N}$  models

$$\exists x \geq \ell (\psi(x) \wedge \Pr(A_{g,f}(x)) < 1/\ell).$$

This is equivalent to (c).

“(c) $\Rightarrow$ (a)”. Let  $f$  be a function as in (c), i.e., for all  $g \in L_{PV}$  and all  $\varepsilon > 0$  there are infinitely many  $m \in \mathbb{N}$  such that  $\Pr(\text{neg}(f(U_m)) \in \text{SAT}) = 1$  and

$$\Pr(|g(U_m)| \leq |f(U_m)| \text{ and } (g(U_m), \text{neg}(f(U_m))) \in \text{Sat}) \leq \varepsilon.$$

Then  $f \in L_{PV}$  and we claim that there is a nonstandard  $n \in M$  such that  $f^M \upharpoonright \{0, 1\}^n$  is a mad pseudo proof system in  $K(F_{PV}^n)$ . We can assume that  $\forall y \text{Fml}(f(y))$  holds in  $\mathbb{N}$  and hence in  $M$ . So to get the madness property (7) it suffices to get  $M \models \psi(n)$  for the formula  $\psi(x)$  defined in (8).

Now, for  $g \in L_{PV}$  let  $B_{g,f} : \mathbb{N} \rightarrow \mathbb{N}$  map  $m \in \mathbb{N}$  to

$$B_{g,f}(m) := \lceil \{ \omega \in \{0, 1\}^m \mid \mathbb{N} \models |g(\omega)| \leq |f(\omega)| \wedge \neg \text{Sat}(g(\omega), f(\omega)) \} \rceil.$$

This function is definable in the standard model  $\mathbb{N}$ , so extends to a function  $B_{g,f}^M$  on  $M$ . We have for all  $n \in M$ :

$$\perp B_{g,f}^M(n) \perp := \{ \omega \in \{0, 1\}^n \mid M \models |g(\omega)| \leq |f(\omega)| \wedge \neg \text{Sat}(g(\omega), f(\omega)) \} \in \mathcal{A}_n.$$

For  $g \in L_{PV}$  and standard  $\ell > 0$  let  $\varphi_{g,f,\ell}(x)$  be the formula

$$x \geq \ell \wedge \psi(x) \wedge \text{Pr}(B_{g,f}(x)) < 1/\ell.$$

Given finitely many such formulas  $\varphi_{g_0,f,\ell_0}, \dots, \varphi_{g_k,f,\ell_k}$  set  $\ell := \max_{i \leq k} \ell_i$  and let  $g \in L_{PV}$  be computed by the following polynomial time algorithm: on input  $\omega$ , compute  $f(\omega), g_0(\omega), \dots, g_k(\omega)$ ; if there is  $i \leq k$  such that  $|g_i(\omega)| \leq |f(\omega)|$  and  $(g_i(\omega), f(\omega)) \notin \text{Sat}$ , then output such  $g_i(\omega)$  (say for the least such  $i \leq k$ ); otherwise output 0.

Then we have for all  $m \in \mathbb{N}$  that  $\bigcup_{i \leq k} \perp B_{g_i,f}(m) \perp \subseteq \perp B_{g,f}(m) \perp$  and thus for all  $i \leq k$ :

$$\text{Pr}(B_{g_i,f}(m)) \leq \text{Pr}(B_{g,f}(m)). \quad (9)$$

Now, by assumption, for every  $\ell$ , there are infinitely many  $m \in \mathbb{N}$  such that  $\mathbb{N} \models \psi(m)$  and  $\text{Pr}(B_{g,f}(m)) < 1/\ell$ . So there is such  $m \geq \ell$ . By choice of  $\ell$  and (9) we get  $\text{Pr}(B_{g_i,f}(m)) < 1/\ell \leq 1/\ell_i$  for all  $i \leq k$ . That is,  $m \geq \ell \geq \ell_i$  satisfies  $\varphi_{g_i,\ell_i}(x)$  for all  $i \leq k$ .

Hence, any finitely many of the formulas  $\varphi_{g,f,\ell}(x)$  for  $g \in L_{PV}$  and standard  $\ell > 0$  are jointly satisfiable in  $\mathbb{N}$  and hence in  $M$ . In other words, the family of subsets of  $M$  defined by these formulas has the finite intersection property. Since  $M$  is  $\aleph_1$ -saturated, there is  $n \in M$  satisfying all these formulas. This  $n$  is nonstandard and satisfies  $M \models \psi(n)$ , hence  $f$  is mad in  $K(F_{PV}^n)$ . Furthermore,  $\text{Pr}^M(B_{g,f}(n)) < 1/\ell$  for all standard  $\ell > 0$  and all  $g \in L_{PV}$ . Hence,  $\perp B_{g,f}(n) \perp \in \text{Inf}_n$  for all  $g \in L_{PV}$ . This is Lemma 3.5 (b), so  $f^M \upharpoonright \{0, 1\}^n$  is a pseudo proof system in  $K(F_{PV}^n)$ .

Finally, we show that statement (b) holds assuming  $\text{NE} \cap \text{coNE} \not\subseteq \text{E}$  and pseudorandom generators exist. By Proposition 2.7, the first assumption gives a polynomial time computable hard sequences for  $\mathbb{L}$ . By Proposition 2.11, the second assumption gives an invertibly samplable probably hard sequence for  $\mathbb{L}$ . By optimality of  $\mathbb{L}$ , this sequence is probably hard for every SAT-solver.

**Acknowledgments** This paper reports results of the first author's Master's thesis [28]. The first author is supported by the Austrian Science Fund (FWF) under project number P25207 and Y698. The second author is supported by the Austrian Science Fund (FWF) under Project P28699 and partially supported by the European Research Council (ERC) under the European Union's Horizon 2020 research programme (grant agreement ERC-2014-CoG 648276 AUTAR).

We thank Yijia Chen, Jan Krajíček and Ján Pich for valuable comments during early stages of this work. We thank Edward Hirsch for pointing out reference [38]. Finally, we are grateful to the reviewers. Thanks to their comments we could substantially improve the presentation of our work.

## References

- [1] S. Arora and B. Barak, Computational Complexity, A Modern Approach (Cambridge University Press, 2009).
- [2] A. Atserias, Distinguishing SAT from polynomial-size circuits, through black-box queries, in: Proceedings of the 21st Annual IEEE Conference on Computational Complexity, held in Prague, Czech Republic, July 16–20, 2006 (IEEE Computer Society, 2006), pp. 88–95.
- [3] A. Bogdanov, K. Talwar, and A. Wan, Hard instances for satisfiability and quasi-one-way functions, in: Innovations in Computer Science, ICS 2010, Tsinghua University, Beijing, China, January 5–7, 2010. Proceedings, edited by A. C.-C. Yao (Tsinghua University Press, 2010), pp. 290–300.
- [4] A. Bogdanov and L. Trevisan, Average-case complexity, Found. Trends Theor. Comput. Sci. 2(1), 1–106 (2006).
- [5] R. Buss, Bounded arithmetic, cryptography and complexity, Theoria 63, 147–167 (1997).

- [6] B. Chapman and R. Williams, The circuit-input game, natural proofs, and testing circuits with data, in: Proceedings of the 2015 Conference on Innovations in Theoretical Computer Science, ITCS 2015, Rehovot, Israel, January 11–13, 2015, edited by T. Roughgarden (ACM Press, 2015), pp. 263–270.
- [7] Y. Chen and J. Flum, On optimal inverters, *Bull. Symb. Log.* **20**(1), 1–23 (2014).
- [8] Y. Chen, J. Flum, and M. Müller, Hard instances of algorithms and proof systems, *ACM Trans. Comput. Theory* **6**(2), article no. 7 (2014).
- [9] A. Cook, Feasibly constructive proofs and the propositional calculus, in: Proceedings of the 7th Annual ACM Symposium on Theory of Computing, May 5–7, 1975, Albuquerque, New Mexico, USA, edited by W. C. Rounds, N. Martin, J. W. Carlyle, M. A. Harrison (ACM Press, 1975), pp. 83–97.
- [10] A. Cook and D. Mitchell, Finding hard instances of the satisfiability problem: a survey, in: Satisfiability Problem: Theory and Applications, Proceedings of the DIMAC Workshop held March 11–13, 1996, edited by D. Du, J. Gu, and P. M. Pardalos, DIMACS Series in Discrete Mathematics and Theoretical Computer Science Vol. 35 (American Mathematical Society, 1997), pp. 1–17.
- [11] A. Cook and R. A. Reckhow, The relative efficiency of propositional proof systems, *J. Symb. Log.* **44**(1), 36–50 (1979).
- [12] M. Garlík, Construction of models of bounded arithmetic by restricted reduced powers, *Arch. Math. Log.* **55**(5), 625–648 (2016).
- [13] D. Gutfreund, R. Shaltiel, and A. Ta-Shma, If **NP** languages are hard on the worst-case, then it is easy to find their hard instances, *Comput. Complexity* **16**(4), 412–441 (2007).
- [14] P. Hájek and P. Pudlák, *Metamathematics of First-Order Arithmetic*, 2nd ed., Perspectives in Mathematical Logic Vol. 3 (Springer, 1998).
- [15] A. Hirsch, D. Itsykson, I. Monakhov, and A. Smal, On optimal heuristic randomized semidecision procedures, with application to proof complexity and cryptography, *Theory Comput. Syst.* **51**, 179–195 (2012).
- [16] J. Hirschfeld, Models of arithmetic and recursive functions, *Isr. J. Math.* **20**, 111–126 (1975).
- [17] T. Jech, *Set Theory, The Third Millenium Edition, Revised and Expanded*, Springer Monographs in Mathematics (Springer, 2002).
- [18] V. Kabanets, Easiness assumptions and hardness tests: trading time for zero error, *J. Comput. Syst. Sci.* **63**(2), 236–252 (2001).
- [19] J. Keisler, *Foundations of Infinitesimal Calculus* (Prindle Weber & Schmidt, 1976).
- [20] J. Keisler, Randomizing a model, *Adv. Math.* **143**, 124–158 (1999).
- [21] R. Klivans and D. van Melkebeek, Graph nonisomorphism has subexponential size proofs unless the polynomial-time hierarchy collapses, *SIAM J. Comput.* **31**(5), 1501–1526 (2002).
- [22] J. Krajíček, Extensions of models of PV, in: Logic Colloquium '95, Proceedings of the Annual European Summer Meeting of the Association of Symbolic Logic, held in Haifa, Israel, August 9–18, 1995, edited by J. A. Makowsky and E. V. Ravve, Lecture Notes in Logic Vol. 11 (Springer, 1998), pp. 104–114.
- [23] J. Krajíček, Forcing with Random variables and Proof Complexity, London Mathematical Society Lecture Note Series Vol. 382 (Cambridge University Press, 2011).
- [24] J. Krajíček, A note on SAT algorithms and proof complexity, *Inform. Process. Lett.* **112**, 490–493 (2012).
- [25] S. Kripke and S. Kochen, Nonstandard models of Peano arithmetic, in: Logic and Algorithmic, Proceedings of An International Symposium Held in Honor of Ernst Specker in Zürich, February 5–11, 1980, edited by H. Läuchli, L'Enseignement Mathématique Vol. 30 (University of Geneva, 1982), pp. 277–295.
- [26] L. Levin, Universal sequential search problems, *Probl. Inf. Transm.* **9**(3), 265–266 (1973).
- [27] J. Lipton and N. E. Young, Simple strategies for large zero-sum games with applications to complexity theory, in: Proceedings of the Twenty-Sixth Annual ACM Symposium on Theory of Computing, 23–25 May 1994, Montréal, Québec, Canada, edited by F. T. Leighton, M. T. Goodrich (ACM Press, 1994), pp. 734–740.
- [28] J. Maly, Jan Krajíček's Forcing Construction and Pseudo Proof Systems, Master's thesis (Universität Wien, 2016).
- [29] R. Mansfield, The theory of Boolean ultrapowers, *Ann. Math. Log.* **2**(3), 297–323 (1971).
- [30] G. McLaughlin, Sub-arithmetical ultrapowers: a survey, *Ann. Pure Appl. Log.* **49**(2), 143–191 (1990).
- [31] J. Pich, Circuit lower bounds in bounded arithmetics, *Ann. Pure Appl. Log.* **166**(1), 29–45 (2015).
- [32] K. Potthoff, Boolean ultrapowers, *Arch. Math. Log. Grundlagenforschung* **16**, 37–48 (1974).
- [33] P. Pudlák, Randomness, pseudorandomness and models of arithmetic, in: New Studies in Weak Arithmetics, edited by P. Cégielski, C. Cornaros, and C. Dimitracopoulos, CSLI Lecture Notes Vol. 211 (CSLI Publications, 2013), pp. 199–216.
- [34] H. Rasiowa and R. Sikorski, Algebraic treatment of the notion of satisfiability, *Fund. Math.* **40**(1), 62–95 (1953).
- [35] D. Scott, A proof of the independence of the continuum hypothesis, *Math. Syst. Theory* **1**(2), 89–111 (1967).
- [36] L. Stockmeyer, The Complexity of Decision Problems in Automata Theory, Ph.D. thesis (Massachusetts Institute of Technology, 1974).
- [37] N. Vereshchagin, An improving on Gutfreund, Shaltiel, and Ta-Shma's paper "If **NP** languages are hard on the worst-case, then it is easy to find their hard instances", in: Computer Science, Theory and Applications, Proceedings of the 8th International Computer Science Symposium in Russia, CSR 2013, held in Ekaterinburg, Russia, June 25–29, 2013, edited by A. A. Bulatov and A. M. Shur, Lecture Notes in Computer Science Vol. 7913 (Springer, 2013), pp. 203–211.
- [38] N. Vereshchagin, Encoding invariance in average case complexity, *Theory Comput. Syst.* **54**(2), 305–317 (2014).