

A Generic Time-Space Lower Bound for Proof Complexity

Scott Diehl*, Dieter van Melkebeek**, and Ryan Williams***

Abstract. We show that for all reals c and d such that $c^2d < 4$ there exists a positive real e such that tautologies cannot be decided by both a nondeterministic algorithm that runs in time n^c , and a nondeterministic algorithm that runs in time n^d and space n^e . In particular, for every $d < \sqrt[3]{4}$ there exists a positive e such that tautologies cannot be decided by a nondeterministic algorithm that runs in time n^d and space n^e .

1 Introduction

Proof complexity studies the NP versus coNP problem — whether tautologies can be recognized efficiently by nondeterministic machines. Typical results in proof complexity deal with specific types of nondeterministic machines that implement well-known proof systems, such as resolution. They establish strong (superpolynomial or even exponential) lower bounds for the size of any proof of certain families of tautologies within that system, and thus for the running time of the corresponding nondeterministic machine deciding tautologies. We refer to [1] for a survey of such results.

Another, more generic, approach to the NP versus coNP problem follows along the lines of the recent time-space lower bounds for satisfiability on deterministic machines [7]. Similar arguments as in the deterministic setting yield somewhat weaker lower bounds for satisfiability on conondeterministic machines, or equivalently, for tautologies on nondeterministic machines. Those results show that no nondeterministic algorithm can decide tautologies in time n^d and space n^e for certain nontrivial combinations of d and e . The lower bounds obtained are very robust with respect to the model of computation, and apply to *any* proof system. However, the arguments only work in the sublinear space range ($e < 1$) and polynomial time range (constant d). For example, Fortnow [2] established a slightly superlinear time lower bound in the case of constant $e < 1$, and Fortnow and Van Melkebeek [4, 3] showed a time lower bound of n^d for any $d < \sqrt{2}$ in the case of subpolynomial space bounds ($e = o(1)$).

In this paper we build on these generic techniques and boost the exponent in the time lower bound for subpolynomial-space nondeterministic algorithms recognizing tautologies from $\sqrt{2} \approx 1.414$ to $\sqrt[3]{4} \approx 1.587$. More precisely, we obtain the following result.

* University of Wisconsin-Madison, sfdiehl@cs.wisc.edu.

** University of Wisconsin-Madison, dieter@cs.wisc.edu.

*** Carnegie Mellon University, ryanw@cs.cmu.edu.

Theorem 1. *For every real $d < \sqrt[3]{4}$ there exists a positive real e such that tautologies cannot be decided by nondeterministic algorithms running in time n^d and space n^e .*

Although raising the exponent from $\sqrt{2}$ to $\sqrt[3]{4}$ may seem like a step that could soon be further improved upon, there is some experimental evidence that $\sqrt[3]{4}$ is in fact the best exponent that can be obtained within the framework of the recent time-space lower bounds. The framework can be formalized such that the exponent of the best time lower bound in the subpolynomial-space range can be computed as the limit of an exhaustive search process. A large automated search did not discover any proofs that lead to an exponent better than $\sqrt[3]{4}$. See [9] for more details.

The earlier result of Fortnow and Van Melkebeek [4, 3] can be refined to rule out either nondeterministic algorithms solving tautologies in time n^c (regardless of space) *or* nondeterministic algorithms solving tautologies in simultaneous time n^d and space n^e for certain combinations of c , d , and e . The precise relationship between the parameters in [4, 3] is that for every c and d such that $(c^2 - 1)d < c$, there is a positive e such that the statement holds. In particular, tautologies cannot have both a nondeterministic algorithm that runs in time $n^{1+o(1)}$ and a nondeterministic algorithm that runs in logarithmic space [2]. Correspondingly, our argument yields the following refinement.

Theorem 2. *For all reals c and d such that $c^2d < 4$, there exists a positive real e such that tautologies cannot be solved by both*

- (i) *a nondeterministic algorithm that runs in time n^c and*
- (ii) *a nondeterministic algorithm that runs in time n^d and space n^e .*

In order to compare the condition $c^2d < 4$ in Theorem 2 with the corresponding condition $(c^2 - 1)d < c$ in [4, 3], we include a plot of those bounds in Figure 1. Note that the interesting range of parameters satisfies $d \geq c \geq 1$. This is because an algorithm of type (ii) is a special case of an algorithm of type (i) for $d \leq c$, and that an algorithm of type (i) with $c < 1$ can be ruled out unconditionally by simple diagonalization. The condition due to this paper, $c^2d < 4$, is less restrictive for values of d not too much larger than c . Thus, Theorem 2 gives a better lower bound in this range. In particular, for $c = d$, our condition requires $d < \sqrt[3]{4} \approx 1.587$, whereas that of [4, 3] requires $d < \sqrt{2} \approx 1.414$; this setting yields the improvement stated in Theorem 1.

Our main technical contribution builds upon the reasoning underlying the $(c^2 - 1)d < c$ condition of [4, 3]. Briefly, the key ingredient in the latter is an inductive argument that builds a sequence of faster and faster nondeterministic algorithms for coNP. Our new approach allows us to harness the algorithm provided by the inductive hypothesis twice in each inductive step, whereas earlier approaches could only do so once per step. This yields better performance in an interesting range of c and d . Due to the double application of the inductive hypothesis, the resulting recurrence relation for the exponents in the algorithms' running times becomes of degree two (rather than one as before) and has non-constant coefficients, but we can still handle it analytically.

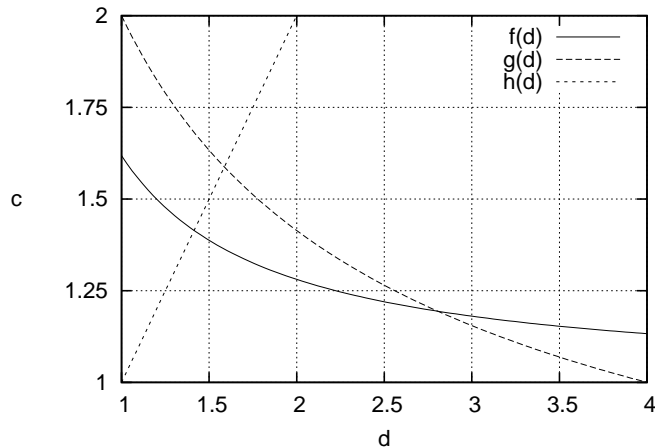


Fig. 1. Tradeoff Curves for Tautologies. Tautologies cannot have both (i) a non-deterministic algorithm that runs in time n^c and (ii) a nondeterministic algorithm that runs in time n^d and space $n^{o(1)}$. The function $f(d)$ solves for the best bound on c due to the prior results of [4, 3], and $g(d)$ does the same for Theorem 2; the function $h(d)$ marks the case $c = d$, illustrating the bound in Theorem 1.

2 Preliminaries

In this section we describe some definitions, conventions, and basic techniques that we use throughout the paper.

2.1 Notation

Much of our notation is standard, although we adopt some nonstandard abbreviations for conciseness. For functions t and s we denote by $\text{NT}(t)$ the class of languages recognized by nondeterministic machines that run in time $O(t)$, and by $\text{NTS}(t, s)$ those recognized by nondeterministic machines that run in simultaneous time $O(t)$ and space $O(s)$. We prefix “co” to either to represent their complementary class. We often use the same notation to refer to classes of machines rather than classes of languages.

Our results are robust with respect to the choice of machine model underlying our complexity classes; for concreteness, we use the random-access machine model as described in [5]. We omit constructibility considerations for the bounds t and s in this paper as our final results apply to polynomial bounds which satisfy all the constructibility properties needed.

Recall that a space-bounded nondeterministic machine does not have two-way access to its guess bits unless it explicitly writes them down on its worktape at the expense of space. It is often important for us to take a finer-grained view of such computations to separate out the resources required to write down a

nondeterministic guess string from those required to verify that the guess is correct. To this end, we adopt the following notation.

Definition 1. *Given a complexity class \mathcal{C} and a function f , we define the class $\exists^f \mathcal{C}$ to be the set of languages that can be described as*

$$\{x | \exists y \in \{0, 1\}^{O(f(|x|))} P(x, y)\},$$

where P is a predicate accepting a language in the class \mathcal{C} when its complexity is measured in terms of $|x|$ (not $|x| + |y|$). We analogously define $\forall^f \mathcal{C}$.

2.2 Tautologies versus Conondeterministic Linear Time

All the known time-space lower bounds for satisfiability or tautologies to date hinge on the tight connection between the tautologies problem and the class of languages recognized by conondeterministic linear-time machines, $\text{coNT}(n)$. The Cook-Levin Theorem, the seminal result showing that satisfiability is NP-complete, can be interpreted as saying that satisfiability captures the time complexity of all of NP up to polynomial factors; the complement of this statement applies to the tautologies problem and coNP . Stronger versions have been formulated for various machine models, showing that tautologies captures the simultaneous time *and* space complexity of conondeterministic linear time on nondeterministic machines up to polylogarithmic factors. As a consequence, time-space lower bounds for $\text{coNT}(n)$ on nondeterministic machines transfer to tautologies with little loss in parameters. In particular we use the following result; we refer to [7] for an elementary and model independent proof.

Lemma 1. *For positive reals d and e , if*

$$\text{coNT}(n) \not\subseteq \text{NTS}(n^d, n^e),$$

then for any reals $d' < d$ and $e' < e$,

$$\text{Tautologies} \not\subseteq \text{NTS}(n^{d'}, n^{e'}).$$

Since a lower bound for $\text{coNT}(n)$ yields essentially the same lower bound for tautologies, we shift our focus to proving lower bounds for the former.

2.3 Indirect Diagonalization

Our proofs follow the paradigm of indirect diagonalization. The paradigm works by contradiction, i.e., we begin by assuming that the desired lower bound does not hold. In the case of Theorem 2 we assume that

$$\text{coNT}(n) \subseteq \text{NT}(n^c) \cap \text{NTS}(n^d, n^e). \tag{1}$$

We then use this unlikely assumption to derive a series of more and more unlikely inclusions of complexity classes. The argument concludes when we derive an inclusion so unlikely that it contradicts a known diagonalization result.

Most of the challenge in formulating an indirect diagonalization argument lies in deriving new inclusions from the assumption (1). The main two tools we use towards this end go in opposite directions:

- (a) Speed up nondeterministic space-bounded computations by adding alternations, and
- (b) Eliminate these alternations at a moderate increase in running time via the assumption (1).

To envision the utility of these items, notice that the assumption allows the simulation of a conondeterministic machine by a space-bounded nondeterministic machine. Item (a) allows us to simulate the latter machine by an alternating machine that runs in less time. Item (b) eliminates the alternations from this simulation, increasing the running time modestly. In this way, we end up back at a nondeterministic computation, so that overall we have derived a simulation of a conondeterministic machine by a nondeterministic one. The complexity class inclusion that this simulation yields is a complementation of the form

$$\text{coNT}(t) \subseteq \text{NT}(f(t)), \tag{2}$$

where we seek to make the function f as small as possible by carefully compounding applications of (a) and (b).

In fact, we know how to rule out inclusions of the type (2) for small functions f , say $f(t) = t^{1-\epsilon}$, by a folklore diagonalization argument. This supplies us with the aforementioned result with which we ultimately derive a contradiction.

Lemma 2. *Let a and b be positive reals such that $a < b$, then*

$$\text{coNT}(n^b) \not\subseteq \text{NT}(n^a).$$

Let us discuss how to achieve items (a) and (b). Item (a) is filled in by the divide-and-conquer strategy that underlies Savitch's Theorem [6]. Briefly, the idea is to divide the computation tableau of a space-bounded nondeterministic machine M into b blocks. Observe that M accepts x in time t if and only if there are $b - 1$ configurations C_1, C_2, \dots, C_{b-1} at the boundaries of these blocks such that for every block i , $1 \leq i \leq b$, the configuration at the beginning of that block, C_{i-1} , can reach the configuration at the end of that block, C_i , in t/b steps, where C_0 is the initial configuration and C_b is the accepting configuration. This condition is implemented on an alternating machine to realize a speedup of M : First existentially guess $b - 1$ configurations of M , universally guess a block number i , and conclude by deciding if C_{i-1} reaches C_i via a simulation of M for t/b steps. Thus, we can derive that:

$$\text{NTS}(t, s) \subseteq \exists^{bs} \forall^{\log b} \text{NTS}(t/b, s). \tag{3}$$

The above simulation runs in overall time $O(bs + t/b)$. Choosing $b = O(\sqrt{t/s})$ optimizes this running time to $O(\sqrt{ts})$. However, minimizing the overall running time of (3) produces suboptimal results in our arguments. Instead, we apply (3) for an unspecified b and choose the optimal value after all of our derivations.

We point out one important fact about the simulation underlying (3): The final phase of this simulation, that of simulating M for t/b steps, does not need access to all of the configurations guessed during the initial existential phase —

it only reads the description of two configurations, C_{i-1} and C_i , in addition to the original input x . Thus, the input size of the final stage is $O(n+s)$ as opposed to $O(n+bs)$ as the complexity-class inclusion of (3) suggests in general. This fact has a subtle but key impact on our analysis in Section 3.

We now turn to item (b), that of eliminating the alternations introduced by (3). In general, eliminating alternations comes at an exponential cost. However, in our case we are armed with the indirect diagonalization assumption (1). The assumption that $\text{coNT}(n) \subseteq \text{NT}(n^c)$ allows us to eliminate an alternation at the cost of raising the running time to the power of c . Alternatively, the assumption that $\text{coNT}(n) \subseteq \text{NTS}(n^d, n^e)$ allows us to eliminate an alternation at the cost of raising the running time to the power of d while at the same time maintaining the space restriction of $O(n^e)$ on the final stage. We use both of these techniques in our analysis.

We now have all the tools we need to carry out our indirect diagonalization argument to prove Theorem 2.

3 Proof of the Lower Bound

We begin with a brief discussion of the techniques required to prove the condition $(c^2 - 1)d < c$ of [4, 3]. We then show how to build on these techniques to arrive at our new condition $c^2d < 4$.

The relevant technical lemma from [4, 3] can be thought of as trading space for time within NP under the indirect diagonalization assumption (1). More precisely, it tries to establish

$$\text{NTS}(t, s) \subseteq \text{NT}(g(t, s)) \tag{4}$$

for the smallest possible functions g , with the hope that $g(t, s) \ll t$. In particular, for subpolynomial space bounds ($s = t^{o(1)}$) and sufficiently large polynomial t , [4, 3] achieves $g = t^{c-1/c+o(1)}$,

$$\text{NTS}(t, t^{o(1)}) \subseteq \text{NT}(t^{c-1/c+o(1)}), \tag{5}$$

which is smaller than t when $c < \phi \approx 1.618$.

As an example of the utility of the space-for-time statement represented by (5), let us sketch the $n^{\sqrt{2}-o(1)}$ lower bound of [4, 3] for subpolynomial-space nondeterministic algorithms solving tautologies. We assume, by way of contradiction, that

$$\text{coNT}(n) \subseteq \text{NTS}(n^c, n^{o(1)}). \tag{6}$$

Then, for sufficiently large polynomials t , we have that:

$$\begin{aligned} \text{coNT}(t) &\subseteq \text{NTS}(t^c, t^{o(1)}) \quad [\text{by assumption (6)}] \\ &\subseteq \text{NT}(t^{c^2-1+o(1)}) \quad [\text{by trading space for time using (5)}]. \end{aligned}$$

This is a contradiction with Lemma 2 when $c < \sqrt{2}$, yielding the desired lower bound.

The space-for-time inclusion (5) is shown by an inductive argument that derives statements of the type (4) for a sequence of smaller and smaller running times $g = g_\ell$, $\ell = 1, 2, \dots$. The idea can be summarized as follows: We start with a space-bounded nondeterministic machine and apply the speedup (3):

$$\text{NTS}(t, s) \subseteq \underbrace{\exists^{bs} \forall^{\log b} \underbrace{\text{NTS}(t/b, s)}_{(*)}}_{(**)}. \quad (7)$$

We then use the inductive hypothesis to trade the space bound of the final stage (*) of this Σ_3 -simulation for time:

$$\text{NTS}(t, s) \subseteq \exists^{bs} \forall^{\log b} \text{NT}(g_{\ell-1}(t/b, s)).$$

We conclude the inductive argument by using the assumption (6) to eliminate the two alternations in this simulation, ending up with another statement of the form

$$\text{NTS}(t, s) \subseteq \text{NT}(g_\ell(t, s)).$$

Notice that the above inductive argument does not rely on the space bound in (6); the weaker assumption that $\text{coNT}(n) \subseteq \text{NT}(n^c)$ is enough to eliminate the alternations introduced by the speedup. Our new argument does exploit the fact that when we transform (*) using the assumption (6), we eliminate an alternation *and* re-introduce a space-bound. This allows us to apply the inductive hypothesis for a *second time* and trade the space bound for a speedup in time once more. This way, we hope to eliminate the alternation in (**) more efficiently than before, yielding a smaller g_ℓ after completing the argument.

Some steps of our new argument exploit the space bound while others do not. In the analysis we allow for different parameters in those two assumptions, say we assume that

$$\text{coNT}(n) \subseteq \text{NTS}(n^c) \cap \text{NTS}(n^d, n^{o(1)}),$$

where $d \geq c \geq 1$. The success of our approach to eliminate the alternation in (**) now depends on how large d is compared with c : If d is not too large compared to c , then the increased cost of complementing via the space-bounded assumption is counteracted by the benefit of trading this space bound for time. That our approach works better in this range of c and d makes plausible the behavior illustrated in Figure 1.

Two key ingredients that allow the above idea to yield a quantitative improvement for certain values of c and d are (i) that the conondeterministic guess at the beginning of stage (**) is only over $\log b$ bits and (ii) the fact mentioned in Section 2 that (*) has input size only $O(n + s)$. Because of (i), the running time of (**) is dominated by that of (*), allowing us to reduce the cost of simulating (**) without an alternation by reducing the cost of simulating (*) in coNT . Item (ii) is important for the latter task because the effective input size for the computation (*) is much smaller than the $O(n + bs)$ bits taken by (**);

in particular, it does not increase with b . This allows the use of larger block numbers b to achieve greater speedups while maintaining that the final stage runs in time at least linear in its input. The latter behavior is crucial in allowing alternation removal at the expected cost — raising the running time to the power of c or d — because we can pad the indirect diagonalization assumption (1) up (to superlinear time) but not down (to sublinear time).

Now that we have sketched the intuition and key ingredients, we proceed with the actual argument. The following lemma formalizes the inductive process of speeding up nondeterministic space-bounded computations on space-unbounded nondeterministic machines.

Lemma 3. *If*

$$\text{coNT}(n) \subseteq \text{NT}(n^c) \cap \text{NTS}(n^d, n^e)$$

for some reals c, d , and e then for every nonnegative integer ℓ , time function t , and space function $s \leq t$,

$$\text{NTS}(t, s) \subseteq \text{NT}((ts^\ell)^{\gamma_\ell} + (n+s)^{a_\ell}),$$

where $\gamma_0 = 1$, $a_0 = 1$, and γ_ℓ and a_ℓ are defined recursively for $\ell > 0$ as follows: Let

$$\mu_\ell = \max(\gamma_\ell(d + e\ell), ea_\ell), \quad (8)$$

then

$$\gamma_{\ell+1} = c\gamma_\ell\mu_\ell / (1 + \gamma_\ell\mu_\ell), \quad (9)$$

and

$$a_{\ell+1} = ca_\ell \cdot \max(1, \mu_\ell). \quad (10)$$

Proof. The proof is by induction on ℓ . The base case $\ell = 0$ is trivial. To argue the inductive step, $\ell \rightarrow \ell + 1$, we consider a nondeterministic machine M running in time t and space s and construct a simulation with the goal of achieving a speedup at the cost of sacrificing the space bound. We begin by simulating M in the third level of the polynomial-time hierarchy via the speedup (3) using $b > 0$ blocks (to be determined later); this simulation is in

$$\exists^{bs} \forall^{\log t} \underbrace{\text{NTS}(t/b, s)}_{(*)}. \quad (11)$$

We focus on simulating the computation of $(*)$ as described above. Recall that the input to $(*)$ consists of the original input x of M as well as two configuration descriptions of size $O(s)$, for a total input size of $O(n + s)$. The inductive hypothesis allows the simulation of $(*)$ in

$$\text{NT}\left(\left(\frac{t}{b}s^\ell\right)^{\gamma_\ell} + (n+s)^{a_\ell}\right). \quad (12)$$

In turn, this simulation can be complemented while simultaneously introducing a space bound via the assumption of the lemma; namely, (12) is in

$$\text{coNTS}\left(\left(\left(\frac{t}{b}s^\ell\right)^{\gamma_\ell} + (n+s)^{a_\ell}\right)^d, \left(\left(\frac{t}{b}s^\ell\right)^{\gamma_\ell} + (n+s)^{a_\ell}\right)^e\right), \quad (13)$$

where here the $(n+s)^{a_\ell}$ term subsumes the $O(n+s)$ term from the input size because $a_\ell \geq 1$. The space bound allows for a simulation via the inductive hypothesis once more, yielding a simulation of $(*)$ in

$$\begin{aligned} \text{coNT} \left(\left(\left(\frac{t}{b} s^\ell \right)^{\gamma_\ell} + (n+s)^{a_\ell} \right)^{\gamma_\ell(d+e\ell)} + (n+s) + \left(\left(\frac{t}{b} s^\ell \right)^{\gamma_\ell} + (n+s)^{a_\ell} \right)^e \right)^{a_\ell} \\ \subseteq \text{coNT} \left(\left(\frac{t}{b} s^\ell \right)^{\gamma_\ell \mu_\ell} + (n+s)^{a_\ell \mu_\ell} + (n+s)^{a_\ell} \right). \end{aligned} \quad (14)$$

Replacing $(*)$ in (11) by (14) eliminates an alternation, lowering the simulation of M to the second level:

$$\exists^{bs} \forall^{\log t} \underbrace{\text{coNT} \left(\left(\frac{t}{b} s^\ell \right)^{\gamma_\ell \mu_\ell} + (n+s)^{a_\ell \mu_\ell} + (n+s)^{a_\ell} \right)}_{(**)} \quad (15)$$

We now complement the conondeterministic computation represented by $(**)$ via the lemma's assumption that $\text{NT}(n) \subseteq \text{coNT}(n^c)$, eliminating one more alternation in the simulation of M . Specifically, since $(**)$ takes input of size $O(n+bs)$, this places the simulation in

$$\begin{aligned} \exists^{bs} \text{NT} \left(\left(\left(\frac{t}{b} s^\ell \right)^{\gamma_\ell \mu_\ell} + (n+s)^{a_\ell \mu_\ell} + (n+s)^{a_\ell} + (bs+n) \right)^c \right) \\ \subseteq \text{NT} \left(\underbrace{\left(\frac{t}{b} s^\ell \right)^{\gamma_\ell \mu_\ell} + (n+s)^{a_\ell \mu_\ell} + (n+s)^{a_\ell}}_{(\searrow)} + \underbrace{bs}_{(\nearrow)} \right)^c, \end{aligned} \quad (16)$$

where the inclusion holds by collapsing the adjacent existential phases (and the time required to guess the $O(bs)$ configuration bits is accounted for by the observation that $c \geq 1$).

Therefore, we have arrived at a simulation that gives rise to an inclusion of $\text{NTS}(t, s)$ in $\text{NT}(\cdot)$; all that remains is to choose b to optimize the running time. Notice that the running time of the simulation in (16) has one term, (\nearrow) , that increases with b and one term, (\searrow) , that decreases with b . The running time is optimized up to a constant factor by choosing b to equate the two terms, resulting in a choice of

$$b^* = \left(\frac{(ts^\ell)^{\gamma_\ell \mu_\ell}}{s} \right)^{1/(1+\gamma_\ell \mu_\ell)}.$$

When this value is at least 1, the running time of the nondeterministic simulation (16) is

$$O \left((ts^{\ell+1})^{c\gamma_\ell \mu_\ell / (1+\gamma_\ell \mu_\ell)} + (n+s)^{ca_\ell \mu_\ell} + (n+s)^{ca_\ell} \right),$$

resulting in the recurrences (9) and (10). If $b^* < 1$, then $b = 1$ is the best we can do; the desired bound still holds since in this case $(\nearrow) + (\searrow) = O(s)$, which is dominated by the $(n+s)^{a_\ell+1}$ term. \square

Under the assumption of Lemma 3, we can further deduce that for a sufficiently large polynomial τ ,

$$\text{coNT}(\tau) \subseteq \text{NTS}(\tau^d, \tau^e) \subseteq \text{NT}(\tau^{(d+e\ell)\gamma_\ell} + \tau^{ea_\ell}) = \text{NT}(\tau^{\mu_\ell}), \quad (17)$$

which is a contradiction with Lemma 2 when $\mu_\ell < 1$. Therefore, the key question is for what values of c , d , and e does μ_ℓ take on a value less than 1. Our analysis focuses on small values of e and shows how such a setting allows us to exhibit the desired behavior in μ_ℓ .

Theorem 3. *For all reals c and d such that $c^2d < 4$ there exists a positive real e such that*

$$\text{coNT}(n) \not\subseteq \text{NT}(n^c) \cap \text{NTS}(n^d, n^e).$$

Proof. The case where either $c < 1$ or $d < 1$ is ruled out by Lemma 2. For $c \geq 1$ and $d \geq 1$, assume (by way of contradiction) that

$$\text{coNT}(n) \subseteq \text{NT}(n^c) \cap \text{NTS}(n^d, n^e)$$

for a value of e to be determined later. As noted above, the theorem's assumption in conjunction with Lemma 3 yields the complementation (17) for any integer $\ell \geq 0$ and sufficiently large polynomial bound τ .

Our goal is now to characterize the behavior of μ_ℓ in terms of c , d , and e . This task is facilitated by focusing on values of e that are small enough to smooth out the complex behavior of μ_ℓ caused by (i) the appearance of the nonconstant term $e\ell$ in the recurrence and (ii) its definition via the maximum of two functions.

We first handle item (i) by introducing a related, nicer sequence by substituting a real β (to be determined) as an upper bound for $e\ell$: Let

$$\mu'_\ell = \max(\gamma'_\ell(d + \beta), ea'_\ell), \quad (18)$$

where $\gamma'_0 = 1$, $a'_0 = 1$ and

$$\gamma'_{\ell+1} = c\gamma'_\ell\mu'_\ell / (1 + \gamma'_\ell\mu'_\ell), \quad (19)$$

and

$$a'_{\ell+1} = ca'_\ell \cdot \max(1, \mu'_\ell). \quad (20)$$

As long as β behaves as intended, i.e., that $e\ell \leq \beta$, we can show by induction that $\gamma_\ell \leq \gamma'_\ell$, $a_\ell \leq a'_\ell$, and $\mu_\ell \leq \mu'_\ell$. Therefore, μ'_ℓ upper bounds μ_ℓ up to a value of ℓ that depends on e , and this ℓ -value becomes large when e is very small. This allows us to use μ'_ℓ as a proxy for μ_ℓ in our analysis.

To smooth out the behavior caused by issue (ii), we point out that the first term in the definition (18) of μ'_ℓ is larger than the second when e is very small. Provided that this is the case, μ'_ℓ equals the sequence ν_ℓ defined as follows:

$$\begin{aligned} \nu_0 &= d + \beta \\ \nu_{\ell+1} &= \nu_\ell^2 c(d + \beta) / ((d + \beta) + \nu_\ell^2). \end{aligned} \quad (21)$$

This delivers a simpler sequence to analyze. Notice that because the underlying transformation

$$\eta \rightarrow \eta^2 c(d + \beta) / ((d + \beta) + \eta^2) \quad (22)$$

is increasing over the positive reals, the sequence ν_ℓ is monotone in this range. It is decreasing if and only if $\nu_1 < \nu_0$, which is equivalent to $(c - 1)(d + \beta) < 1$.

Furthermore, when $c^2(d+\beta) < 4$, the transformation has a unique real fixed point at 0. Since the underlying transformation is also bounded and starts positively, the sequence ν_ℓ must decrease monotonically to 0 in this case.

Therefore, when $c^2d < 4$ we can choose a positive β such that ν_ℓ becomes as small as we want for large ℓ . Provided that β , e , and ℓ satisfy the assumptions required to smooth out items (i) and (ii), this also gives us that μ_ℓ is small. More formally, let ℓ^* be the first value of ℓ such that $\nu_\ell < 1$. Then to satisfy item (i), we require that

$$e\ell^* \leq \beta. \tag{23}$$

For item (ii), we require that the first term in the definition (18) of μ'_ℓ dominates the second up to this point, namely,

$$\gamma'_\ell(d + \beta) \geq ea'_\ell \text{ for all } \ell \leq \ell^*. \tag{24}$$

When all of these conditions are satisfied, we have that

$$\mu_{\ell^*} \leq \mu'_{\ell^*} = \nu_{\ell^*} < 1,$$

and the running time of the conondeterministic simulation represented by (17) for $\ell = \ell^*$ runs in time

$$O(\tau^{\mu_{\ell^*}}) = O(\tau^{\mu'_{\ell^*}}) = O(\tau^{\nu_{\ell^*}}). \tag{25}$$

Therefore, by choosing a small enough positive e to satisfy the finite number of constraints in (23) and (24), we arrive at our goal of exhibiting an exponent cost in the complementation of (17) that is smaller than 1. This is a contradiction, which proves the desired lower bound. \square

The proof of Theorem 3 establishes $c^2d < 4$ as a sufficient condition for our approach to work but it isn't clear that the condition is also necessary. For completeness we include an argument showing that our analysis in the proof is tight — our approach does not work for any setting with $c^2d \geq 4$.

Referring to the notation used in the proof of Theorem 3, our approach works if we can find a value of $\mu_\ell < 1$ so that (17) contradicts Lemma 3. Without loss of generality, we can assume that the space-efficient simulation only uses logarithmic space, in which case the sequence μ_ℓ simplifies to ν_ℓ given by (21), where the term β is negligible. Since $\nu_1 \geq 1$, ν_ℓ has to decrease in order to obtain a contradiction; as this happens when $(c-1)d < 1$, we can rule out success in settings with $c^2d \geq 4$ and $d = 1$, or $c \geq 2$. For $c^2d = 4$ and $d > 1$, the underlying transformation (22) has a unique positive fixed point to which the sequence ν_ℓ converges, namely $cd/2 = \sqrt{d}$, which is less than d . If we let c grow, the double fixed point separates into two positive fixed points that gradually diverge from but remain centered around $cd/2$. As long as the largest of the two fixed points remains less than d , the sequence converges to it in a monotone decreasing way. At the verge where that fixed point reaches d , the sequence remains constant, which means that $(c-1)d = 1$. Beyond that the sequence monotonically increases to the larger fixed point. This argument implies that at all times $\nu_\ell \geq cd/2$. In order to have $\nu_\ell < 1$, we would need $cd < 2$, which is incompatible with our assumption that $c^2d \geq 4$ and $d > 1$.

4 Future Work

The most compelling avenue for future work is to improve the quantitative strength of the lower bounds for tautologies. There is no a priori reason to believe that the exponent of $\sqrt[3]{4}$ could not be improved upon in the near future by making modifications to current arguments. For example, one might think that applying the inductive hypothesis more than twice per step would lead to further progress, but this fails¹. In fact, an automated search exploiting the regularity within current indirect diagonalization practices revealed no evidence of *any* proof doing better than $\sqrt[3]{4}$ [9]. Therefore, we believe that an improvement to our lower bound for tautologies must involve novel ingredients or even a completely new approach to proving lower bounds for NP-complete and related hard problems.

References

1. P. Beame and T. Pitassi. Propositional proof complexity: Past, present, and future. In G. Paun, G. Rozenberg, and A. Salomaa, editors, *Current Trends in Theoretical Computer Science: Entering the 21st Century*, pages 42–70. World Scientific, 2001.
2. L. Fortnow. Time-space tradeoffs for satisfiability. *Journal of Computer and System Sciences*, 60:337–353, 2000.
3. L. Fortnow, R. Lipton, D. van Melkebeek, and A. Viglas. Time-space lower bounds for satisfiability. *Journal of the ACM*, 52:835–865, 2005.
4. L. Fortnow and D. van Melkebeek. Time-space tradeoffs for nondeterministic computation. In *Proceedings of the 15th IEEE Conference on Computational Complexity*, pages 2–13. IEEE, 2000.
5. D. van Melkebeek. Time-space lower bounds for NP-complete problems. In G. Paun, G. Rozenberg, and A. Salomaa, editors, *Current Trends in Theoretical Computer Science*, pages 265–291. World Scientific, 2004.
6. W. Savitch. Relationships between nondeterministic and deterministic tape complexities. *Journal of Computer and System Sciences*, 4:177–192, 1970.
7. D. van Melkebeek. A survey of lower bounds for satisfiability and related problems. *Foundations and Trends in Theoretical Computer Science*, 2:197–303, 2007.
8. R. Williams. Time-space tradeoffs for counting NP solutions modulo integers. In *Proceedings of the 22nd IEEE Conference on Computational Complexity*, pages 70–82. IEEE, 2007.
9. R. Williams. Automated proofs of time lower bounds. Manuscript, 2008.

¹ The reason can be traced back to the fact that the size of the input to $(**)$ in (7) is the full $O(n + bs)$ rather than the smaller $O(n + s)$ taken by $(*)$.