

# Toward a Characterization of Simulation Between Arithmetic Theories

Hunter Monroe

April 2026

## Abstract

We study when a sound arithmetic theory  $\mathcal{S} \supseteq \mathcal{S}_2^1$  with polynomial-time decidable axioms efficiently proves the bounded consistency statements  $Con_{\mathcal{S}+\phi}(n)$  for a true sentence  $\phi$ . Equivalently, we ask when  $\mathcal{S}$ , viewed as a proof system, simulates  $\mathcal{S}+\phi$ . The paper’s two unconditional contributions constrain possible characterizations. First, for finitely axiomatized sequential  $\mathcal{S}$ , if  $EA \vdash Con_{\mathcal{S}} \rightarrow Con_{\mathcal{S}+\phi}$ , then  $\mathcal{S}$  interprets  $\mathcal{S}+\phi$ , implying  $\mathcal{S} \stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}}(p(n)) \rightarrow Con_{\mathcal{S}+\phi}(n)$  for some polynomial  $p$ , and hence  $\mathcal{S} \stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}+\phi}(n)$ . Second, if  $\mathcal{S}$  fails to simulate  $\mathcal{S}+\phi$  for some true  $\phi$ , then for all sufficiently large  $k$  it also fails for  $\phi_{BB}(k)$  asserting the exact value of the  $k$ -state Busy Beaver function. Informally, any argument showing that  $\mathcal{S}$  fails to simulate some  $\mathcal{S}+\phi$  also yields unprovable  $\phi_{BB}(k)$  witnessing the same obstruction. These results suggest that relative consistency strength is a serious candidate for governing when simulation is possible, while leaving open whether it is the correct criterion.

The paper’s central conjectural proposal is that the above sufficient condition is also necessary: if  $EA \not\vdash Con_{\mathcal{S}} \rightarrow Con_{\mathcal{S}+\phi}$ , then for every constant  $c > 0$ ,  $\mathcal{S} \not\stackrel{n^c}{\vdash} Con_{\mathcal{S}+\phi}(n)$ . Under this proposal, hardness follows in canonical cases where  $\phi$  is  $Con_{\mathcal{S}}$  or a Kolmogorov-randomness axiom. The latter yields further conjectural consequences and extensions.

## 1 Introduction

*“...the essence of the big open problems in complexity theory could be logical, rather than combinatorial”* — Pudlák [11]

Let  $\mathcal{S} \supseteq S_2^1$  be a sound arithmetic theory with polynomial-time decidable axioms, and let  $\phi$  be a true sentence independent of  $\mathcal{S}$ . We study when  $\mathcal{S}$  efficiently proves the bounded consistency statements  $Con_{\mathcal{S}+\phi}(n)$ , where  $Con_{\mathcal{S}+\phi}(n)$  asserts that  $\mathcal{S}+\phi$  has no proof of  $0=1$  within  $n$  symbols.

From the perspective of proof complexity, this asks when  $\mathcal{S}$ , viewed as a proof system, simulates  $\mathcal{S}+\phi$ .<sup>1</sup> Independence alone does not prevent simulation: there are true independent sentences  $\phi$  such that  $\mathcal{S}$  polynomial-time interprets  $\mathcal{S}+\phi$ , and hence efficiently proves  $Con_{\mathcal{S}+\phi}(n)$ ; see Pudlák [12, Lemma 3.6]. On the other hand, if no optimal proof system exists, then for some true sentence  $\phi$  one has for all  $c$ ,  $\mathcal{S} \not\vdash^{n^c} Con_{\mathcal{S}+\phi}(n)$ ; see Krajíček and Pudlák [8, Theorem 2.1].<sup>2</sup> Thus the problem is to identify a criterion for  $\phi$  that governs simulation and non-simulation.

We formulate this as follows.

**Characterization Problem.** For every sound theory  $\mathcal{S} \supseteq S_2^1$  and every true sentence  $\phi$  independent of  $\mathcal{S}$ , determine when  $\mathcal{S} \vdash^{n^{O(1)}} Con_{\mathcal{S}+\phi}(n)$ , and when instead  $\mathcal{S} \not\vdash^{n^c} Con_{\mathcal{S}+\phi}(n)$  for every constant  $c$ . Equivalently, determine when  $\mathcal{S}$  simulates  $\mathcal{S}+\phi$  as a proof system.<sup>3</sup>

A benchmark case is  $\phi = Con_{\mathcal{S}}$ ; Pudlák [10, Problem 1] conjectures that  $\mathcal{S} \vdash^{n^c} Con_{\mathcal{S}+Con_{\mathcal{S}}}(n)$  for every  $c$ . This conjecture implies that higher absolute consistency strength is a sufficient condition for non-simulation (Theorem 3.6 below). The present paper asks for a more general organizing principle for arbitrary true independent sentences  $\phi$ , ideally a condition that is both necessary and sufficient.

The paper’s first contribution identifies a new sufficient condition for simulation for finitely axiomatized sequential  $\mathcal{S}$ . If  $EA \vdash Con_{\mathcal{S}} \rightarrow Con_{\mathcal{S}+\phi}$ , then  $\mathcal{S}$  interprets  $\mathcal{S}+\phi$ . By proof translation, there exists a polynomial  $p$  such that  $\mathcal{S} \vdash^{n^{O(1)}} Con_{\mathcal{S}}(p(n)) \rightarrow Con_{\mathcal{S}+\phi}(n)$ . Thus Theorem 3.4 shows that relative consistency ( $EA \vdash Con_{\mathcal{S}} \rightarrow Con_{\mathcal{S}+\phi}$ ) implies simulation ( $\mathcal{S} \vdash^{n^{O(1)}} Con_{\mathcal{S}+\phi}(n)$ ) for  $\mathcal{S}$  finitely axiomatized and sequential. Taking this theorem to be tight leads to the following structural candidate characterization of simulation. The paper develops that proposal and its main consequences. In particular, the case  $\phi = Con_{\mathcal{S}}$  and the Kolmogorov-randomness axioms arise directly as natural instances of the same general obstruction.

---

<sup>1</sup>For background on proof complexity, see Krajíček[7].

<sup>2</sup>We use the notation  $\mathcal{S} \vdash^{n^{O(1)}} \phi(n)$  to mean that there exists a constant  $c$  such that  $\mathcal{S}$  has proofs of  $\phi(n)$  of size at most  $n^c$  for all sufficiently large  $n$ . Hardness statements are written explicitly in the form “for every constant  $c$ ,  $\mathcal{S} \not\vdash^{n^c} \phi(n)$ .”

<sup>3</sup>We limit attention to true  $\phi$  to isolate non-simulation arising from lack of access to additional *correct* information, rather than trivial failure due to unsound extensions, and because this is the regime relevant for applications to tautologies, where  $\phi$  expresses the truth of a family of propositional formulas.

The paper’s second contribution shows that canonical incompleteness phenomena already suffice to witness non-simulation. More precisely, if  $\mathcal{S}$  fails to simulate  $\mathcal{S}+\phi$  for some true sentence  $\phi$ , then for all sufficiently large  $k$  it also fails to simulate  $S_2^1+\phi_{BB}(k)$ , where  $\phi_{BB}(k)$  asserts the exact value of the  $k$ -state Busy Beaver function, namely the maximum halting time of any halting  $k$ -state Turing machine (Theorem 3.10).<sup>4</sup> The change in base theory here is deliberate: the theorem shows that once hardness occurs above  $S_2^1$ , it is already witnessed by a canonical Busy Beaver extension of the weak base theory  $S_2^1$ . This aligns with Aaronson’s observation that, for sufficiently large  $k$ , true Busy Beaver sentences can prove the consistency of arbitrarily strong computably axiomatized theories [1, Proposition 3]. Thus the hardness side of the Characterization Problem already appears among Busy Beaver sentences. Any characterization of simulation should therefore explain these canonical unprovable truths.

Together, these results constrain possible characterizations and suggest two complementary ways of thinking about hardness. First, in line with the Busy Beaver reduction, efficient simulation should be controlled by information already visible to weak arithmetic; this points toward a structural criterion formulated in terms of relative-consistency transfer. Second, canonical sources of unprovable information such as Busy Beaver values and Kolmogorov-random strings suggest a more semantic or information-theoretic obstruction: simulation should fail when it would require access to true information that the base theory cannot itself recover in any way certified in a weak base theory, such as Elementary Arithmetic. Although weak, such a theory can prove true statements about Turing machines that halt.

These considerations motivate a single structural picture of hardness. The main proposal of the paper is **Higher Relative Consistency**, which asserts that efficient simulation fails exactly when weak arithmetic cannot certify that adjoining  $\phi$  preserves consistency. Within this framework, canonical hard axioms such as Busy Beaver values and Kolmogorov-random truths arise as especially natural instances of the same obstruction.

**Higher Relative Consistency (informal).** If  $EA \not\vdash Con_{\mathcal{S}} \rightarrow Con_{\mathcal{S}+\phi}$ , then for every constant  $c$ ,  $\mathcal{S} \not\vdash^{n^c} Con_{\mathcal{S}+\phi}(n)$ .

In other words, failure of relative-consistency should already rule out efficient simulation. Combined with the positive simulation theorem, this yields a natural conjectural picture in which simulation occurs exactly when weak arithmetic can certify that adjoining  $\phi$  preserves consistency.

The remainder of the paper is organized as follows. Section 2 provides preliminaries. Section 3 presents unconditional constraints on any characterization. Section 4 develops heuristic support for the proposed criterion. Section 5 states **Higher Relative Consistency**, derives its main consequences, and formulates *Kolmogorov Hardness* as a distinguished random-axiom instance of the general characterization. Section 6 develops extensions of Kolmogorov

---

<sup>4</sup>See Rado [13].

Hardness. Section 7 discusses provability of **Higher Relative Consistency**. Section 8 concludes.

## 2 Preliminaries

We work with sound arithmetical theories  $\mathcal{S} \supseteq \mathcal{S}_2^1$  whose axioms are decidable in polynomial time. Throughout,  $\phi$  denotes a true arithmetical sentence, and  $\mathcal{S}+\phi$  is the corresponding true extension of  $\mathcal{S}$ . Unless explicitly stated otherwise, implications of the form  $Con_{\mathcal{S}} \rightarrow Con_{\mathcal{S}+\phi}$  are understood as formalized in  $EA$ . This is the level at which the interpretability criteria used later are naturally stated.

We fix a standard arithmetization of syntax. Proofs are encoded as binary strings, and all proof lengths are measured in symbols under this encoding. As usual, any two reasonable codings yield polynomially equivalent proof lengths, so statements of the form  $\mathcal{S} \upharpoonright^{n^{O(1)}} \phi(n)$  are robust under the choice of coding.

For a theory  $\mathcal{S}$ , let  $Con_{\mathcal{S}}(n)$  denote the bounded consistency statement asserting that there is no  $\mathcal{S}$ -proof of  $0=1$  of length at most  $n$ . We write  $\mathcal{S} \upharpoonright^{n^c} \phi(n)$  to mean that, for all sufficiently large  $n$ , the sentence  $\phi(n)$  has an  $\mathcal{S}$ -proof of length at most  $n^c$ . Likewise,  $\mathcal{S} \upharpoonright^{n^{O(1)}} \phi(n)$  means that such proofs exist with polynomially bounded length.

We say that  $\mathcal{S}$  *simulates*  $\mathcal{S}+\phi$  if  $\mathcal{S} \upharpoonright^{n^{O(1)}} Con_{\mathcal{S}+\phi}(n)$ . Thus the Characterization Problem asks when  $\mathcal{S}$  admits polynomial-size proofs of the bounded consistency statements for the true extension  $\mathcal{S}+\phi$ .

We say that  $\phi$  *raises the relative-consistency strength* of  $\mathcal{S}$  if  $EA \not\vdash Con_{\mathcal{S}} \rightarrow Con_{\mathcal{S}+\phi}$ , where  $EA$  is Elementary Arithmetic (Visser [14]). We introduce new terminology, saying that  $\mathcal{S}$  has *feasible relative consistency* for  $\phi$  if there is a polynomial  $p$  such that  $\mathcal{S} \upharpoonright^{n^{O(1)}} Con_{\mathcal{S}}(p(n)) \rightarrow Con_{\mathcal{S}+\phi}(n)$ .

We distinguish throughout between external proof-length bounds and formalized internal implications. The statement  $\mathcal{S} \upharpoonright^{n^{O(1)}} Con_{\mathcal{S}+\phi}(n)$  is an external assertion about the existence of short  $\mathcal{S}$ -proofs, whereas  $EA \vdash Con_{\mathcal{S}} \rightarrow Con_{\mathcal{S}+\phi}$  is an internal relative-consistency implication in a weak base theory. Much of the paper is concerned with when the former should force the latter.

Finally, we use standard notions of interpretability. When additional hypotheses such as finite axiomatizability, sequentiality, or reflexivity are needed, they will be stated explicitly. In the finitely axiomatized sequential setting, interpretability is closely related to relative-consistency implications by the Orey–Hájek and Visser characterizations, and polynomial-time interpretation yields polynomial overhead proof translation. These standard facts underlie the paper’s unconditional upper-bound mechanism.

Whenever an external proof transformation is used to derive polynomial-size  $\mathcal{S}$ -proofs, we will state explicitly whether the transformation is merely true in the standard model, formalizable in  $EA$ , or available with polynomial-size  $\mathcal{S}$ -proofs.

These notions fix the framework for the paper’s central question: whether simulation of  $\mathcal{S}+\phi$  by  $\mathcal{S}$  is governed by relative-consistency transfer, and in particular by the weak-base implication  $EA\vdash Con_{\mathcal{S}}\rightarrow Con_{\mathcal{S}+\phi}$ .

### 3 Unconditional Constraints on Characterizations

This section presents the paper’s unconditional results. The first shows in certain settings that relative consistency implies feasible relative consistency and hence simulation. The second shows that any hard true  $\phi$  can, in a precise sense, be replaced by one of Busy Beaver form. Together, these results isolate both the positive mechanism underlying simulation and a canonical source of hardness, constraining any possible characterization.

The literature shows that simulation is possible in special cases. In particular, polynomial-time interpretability implies efficient provability of the bounded consistency statements  $Con_{\mathcal{S}+\phi}(n)$ ; see Theorem 3.1.<sup>5</sup> By contrast, Pudlák [10] conjectures that this fails already for  $\phi=Con_{\mathcal{S}}$ .<sup>6</sup> Unconditional lower bounds of the form  $\mathcal{S} \not\vdash^{n^{O(1)}} Con_{\mathcal{S}+\phi}(n)$  are not known for sound theories  $\mathcal{S}\supseteq S_2^1$ ; such a result would imply that no optimal proof system exists, and hence  $\mathbf{P}\neq\mathbf{NP}\neq\mathbf{coNP}$ .

The two results proved in this section constrain possible answers to the Characterization Problem in complementary ways. The first isolates a sufficient condition for simulation, formulated in terms of relative consistency rather than polynomial-time interpretability. The second shows that any instance of non-simulation can already be witnessed by  $\phi$  referring to Busy Beavers. Together, these results clarify both the positive mechanism underlying simulation and a canonical source of hardness.

#### 3.1 A Known Sufficient Condition for Simulation

A proof translation argument by Jeřábek, as presented by Pudlák [12, Lemma 3.6], shows that if  $\mathcal{S}$  polynomial-time interprets  $\mathcal{S}+\phi$ , then polynomial-size proofs of  $Con_{\mathcal{S}}(n)$  yield polynomial-size proofs of  $Con_{\mathcal{S}+\phi}(n)$ . Moreover, for a given  $\mathcal{S}$ , there are true sentences  $\phi$  independent of  $\mathcal{S}$  such that  $\mathcal{S}$  polynomial-time interprets  $\mathcal{S}+\phi$ , for example H-Rosser sentences as in Hájek–Pudlák [5, Theorem 4.5(5)]. Thus, simulation is not ruled out for every true independent sentence  $\phi$ —non-simulation requires more than independence.

We present that argument with a theorem statement and proof in a form convenient for this paper’s exposition.

---

<sup>5</sup>Freund and Pakhomov [4] show that such efficient provability can also occur when  $\phi$  expresses slow consistency.

<sup>6</sup>Khaniki [6] relates a weak form of Pudlák’s Conjecture to the existence of a computable jump operator in proof complexity.

**Theorem 3.1** (Jeřábek, via Pudlák [12, Lemma 3.6]) *Suppose  $\mathcal{S} \supseteq S_2^1$  and  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}}(n)$ . If  $\mathcal{S}$  polynomial-time interprets  $\mathcal{S}+\phi$ , then  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ .*

**Proof:** Let  $i$  be a polynomial-time interpretation of  $\mathcal{S}+\phi$  in  $\mathcal{S}$ . By the standard proof-translation argument for interpretations, as presented in Pudlák [12, Lemma 3.6], there is a polynomial  $p$  such that from every  $(\mathcal{S}+\phi)$ -proof  $\pi$  of a sentence  $\psi$  one can compute an  $\mathcal{S}$ -proof of the translated sentence  $\psi^i$  of length at most  $p(|\pi|)$ .

Apply this to  $\psi \equiv 0=1$ . Then any  $(\mathcal{S}+\phi)$ -proof of contradiction of length at most  $n$  yields an  $\mathcal{S}$ -proof of  $(0=1)^i$  of length at most  $p(n)$ . Since  $i$  is an interpretation,  $(0=1)^i$  is a fixed false sentence, and there is a constant-size  $\mathcal{S}$ -proof of  $(0=1)^i \rightarrow 0=1$ . After increasing  $p$  if necessary, it follows that any  $(\mathcal{S}+\phi)$ -proof of contradiction of length at most  $n$  yields an  $\mathcal{S}$ -proof of  $0=1$  of length at most  $p(n)$ .

The proof translation and its polynomial bound are formalizable in  $S_2^1$ , hence in  $\mathcal{S}$ . Therefore  $\mathcal{S}$  proves that if there is no  $\mathcal{S}$ -proof of contradiction of length at most  $p(n)$ , then there is no  $(\mathcal{S}+\phi)$ -proof of contradiction of length at most  $n$ . That is,  $\mathcal{S} \vdash \forall n: \text{Con}_{\mathcal{S}}(p(n)) \rightarrow \text{Con}_{\mathcal{S}+\phi}(n)$ .

By assumption,  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}}(n)$ . Substituting  $p(n)$  for  $n$  gives polynomial-size  $\mathcal{S}$ -proofs of  $\text{Con}_{\mathcal{S}}(p(n))$ . Composing these with this implication yields polynomial-size  $\mathcal{S}$ -proofs of  $\text{Con}_{\mathcal{S}+\phi}(n)$ . Hence  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ . ■

The essential content of the interpretation argument is not merely simulation, but the existence of a uniform feasible relative-consistency implication from  $\mathcal{S}$  to  $\mathcal{S}+\phi$ . All currently known positive mechanisms for simulation pass through such implications.

A central question is the converse: whether every instance of simulation must admit such an internal explanation.

## 3.2 Stronger Results on Simulation

The proof of Theorem 3.1 relies only on the existence of a polynomial  $p$  such that  $\mathcal{S} \vdash \forall n: \text{Con}_{\mathcal{S}}(p(n)) \rightarrow \text{Con}_{\mathcal{S}+\phi}(n)$ . This suggests the following sufficient condition for simulation:  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}}(p(n)) \rightarrow \text{Con}_{\mathcal{S}+\phi}(n)$ . We call this *feasible relative consistency*.

**Theorem 3.2** *Let  $\mathcal{S} \supseteq S_2^1$  be a sound theory with polynomial-time decidable axioms such that  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}}(n)$ . Then the following are equivalent:*

1.  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}}(p(n)) \rightarrow \text{Con}_{\mathcal{S}+\phi}(n)$  for some polynomial  $p$ .
2.  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ .

**Proof:** For (1) $\rightarrow$ (2), combine polynomial-size proofs of  $\text{Con}_{\mathcal{S}}(p(n))$  with polynomial-size proofs of  $\text{Con}_{\mathcal{S}}(p(n)) \rightarrow \text{Con}_{\mathcal{S}+\phi}(n)$ . This yields polynomial-size proofs of  $\text{Con}_{\mathcal{S}+\phi}(n)$ .

For (2) $\rightarrow$ (1), assume  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ . Fix any polynomial  $p$ . From a proof of  $\text{Con}_{\mathcal{S}+\phi}(n)$ , one can derive a proof of  $\text{Con}_{\mathcal{S}}(p(n)) \rightarrow \text{Con}_{\mathcal{S}+\phi}(n)$  with only polynomial overhead by prefixing

a fixed implicational derivation. Since the formula  $Con_{\mathcal{S}}(p(n))$  has size polynomial in  $n$ , the resulting proof length remains polynomially bounded. Hence  $\mathcal{S} \stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}}(p(n)) \rightarrow Con_{\mathcal{S}+\phi}(n)$ .  $\blacksquare$

Accordingly, in what follows we treat simulation itself, namely  $\mathcal{S} \stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}+\phi}(n)$ , as the primary notion. When  $\mathcal{S} \stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}}(n)$ , feasible relative consistency is simply an equivalent reformulation: condition (1) expresses feasible relative consistency, and condition (2) expresses simulation. To state the contrapositive of the above theorem:

**Corollary 3.3** *Let  $\mathcal{S} \supseteq \mathcal{S}_2^1$  be a sound theory with polynomial-time decidable axioms such that  $\mathcal{S} \stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}}(n)$ . Then, for any polynomial  $p$ , the following are equivalent:*

1.  $\mathcal{S} \not\stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}+\phi}(n)$ .
2.  $\mathcal{S} \not\stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}}(p(n)) \rightarrow Con_{\mathcal{S}+\phi}(n)$ .

**Proof:** This follows immediately from Theorem 3.2 by contraposition in both directions.  $\blacksquare$

The interpretation theorem yields the following stronger positive result in the finitely axiomatized sequential setting.

**Theorem 3.4** *Suppose  $\mathcal{S} \supseteq \mathcal{S}_2^1$  is finitely axiomatized and sequential. If  $EA\text{-}Con_{\mathcal{S}} \rightarrow Con_{\mathcal{S}+\phi}$ , then there exists a polynomial  $p$  such that  $\mathcal{S} \stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}}(p(n)) \rightarrow Con_{\mathcal{S}+\phi}(n)$ . In particular,  $\mathcal{S} \stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}+\phi}(n)$ .*

**Proof:** Pudlák [10] shows  $\mathcal{S} \stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}}(n)$ . By the standard interpretability criterion for finitely axiomatized sequential theories, as presented by Visser [14], the hypothesis  $EA\text{-}Con_{\mathcal{S}} \rightarrow Con_{\mathcal{S}+\phi}$  implies that  $\mathcal{S}$  interprets  $\mathcal{S}+\phi$ . Since  $\mathcal{S}$  is finitely axiomatized, so is  $\mathcal{S}+\phi$ , and a fixed interpretation of a finitely axiomatized theory yields a polynomial-time proof translation: the map  $\psi \mapsto \psi^i$  is computable in polynomial time, and each translated axiom has a fixed  $\mathcal{S}$ -proof. Hence,  $\mathcal{S}$  polynomial-time interprets  $\mathcal{S}+\phi$ .

By Theorem 3.1, it follows that  $\mathcal{S} \stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}+\phi}(n)$ . The existence of a polynomial  $p$  such that  $\mathcal{S} \stackrel{n^{O(1)}}{\vdash} Con_{\mathcal{S}}(p(n)) \rightarrow Con_{\mathcal{S}+\phi}(n)$  then follows from Theorem 3.2.  $\blacksquare$

By contrast, Pudlák [10] conjectures that even the benchmark case  $\phi = Con_{\mathcal{S}}$  already yields non-simulation:

**Conjecture 3.5** (*Pudlák's Conjecture*) *For every constant  $c$ ,  $\mathcal{S} \not\stackrel{n^c}{\vdash} Con_{\mathcal{S}+Con_{\mathcal{S}}}(n)$ .*<sup>7</sup>

<sup>7</sup>Another natural candidate for a hard extension is given by the Buss jump [2], that is, a sentence or schema asserting the soundness of the next level of bounded reasoning over  $\mathcal{S}$ .

Under Pudlák’s Conjecture, any true extension that raises absolute consistency (proves  $Con_{\mathcal{S}}$ ) is hard:

**Theorem 3.6** *Assume Pudlák’s Conjecture, and suppose  $\mathcal{S}+\phi\vdash Con_{\mathcal{S}}$ . Then for every constant  $c$ ,  $\mathcal{S}\not\vdash^{\frac{n^c}{7}} Con_{\mathcal{S}+\phi}(n)$ .*

**Proof:** Fix a proof  $\pi$  of  $Con_{\mathcal{S}}$  in  $\mathcal{S}+\phi$ .

There is a uniform proof transformation sending any  $(\mathcal{S}+Con_{\mathcal{S}})$ -proof to an  $(\mathcal{S}+\phi)$ -proof by replacing each use of the extra axiom  $Con_{\mathcal{S}}$  by the fixed derivation  $\pi$ . This increases proof length by at most a constant multiplicative and additive factor. Hence there exist a linear polynomial  $p$  and a constant  $d$  such that  $\mathcal{S}\vdash^{\frac{n^d}{7}} Con_{\mathcal{S}+\phi}(p(n))\rightarrow Con_{\mathcal{S}+Con_{\mathcal{S}}}(n)$ .

Suppose toward a contradiction that  $\mathcal{S}\not\vdash^{\frac{n^{O(1)}}{7}} Con_{\mathcal{S}+\phi}(n)$ . Since  $p$  is linear, substituting  $p(n)$  for  $n$  still gives  $\mathcal{S}\not\vdash^{\frac{n^{O(1)}}{7}} Con_{\mathcal{S}+\phi}(p(n))$ . Composing these proofs with this implication yields  $\mathcal{S}\not\vdash^{\frac{n^{O(1)}}{7}} Con_{\mathcal{S}+Con_{\mathcal{S}}}(n)$ , contrary to Pudlák’s Conjecture.

Therefore for every constant  $c$ ,  $\mathcal{S}\not\vdash^{\frac{n^c}{7}} Con_{\mathcal{S}+\phi}(n)$ . ■

Then there is a gap between Pudlák’s Conjecture and the strongest known sufficient condition for simulation, in Theorem 3.4:

**Corollary 3.7** *Pudlák’s Conjecture is not tight as a converse to Theorem 3.4.*

**Proof:** Let  $\theta$  be any true sentence such that  $\mathcal{S}+Con_{\mathcal{S}}\not\vdash\theta$ , and put  $\phi:=Con_{\mathcal{S}}\wedge\theta$ . Then  $\mathcal{S}+\phi\vdash Con_{\mathcal{S}}$ , so by Theorem 3.6, one has, for every constant  $c$ ,  $\mathcal{S}\not\vdash^{\frac{n^c}{7}} Con_{\mathcal{S}+\phi}(n)$ .

However,  $\phi$  is strictly stronger than  $Con_{\mathcal{S}}$  over  $\mathcal{S}$ . Indeed,  $\mathcal{S}\vdash\phi\rightarrow Con_{\mathcal{S}}$  is immediate. If also  $\mathcal{S}\vdash Con_{\mathcal{S}}\rightarrow\phi$ , then  $\mathcal{S}+Con_{\mathcal{S}}\vdash\phi$ , hence  $\mathcal{S}+Con_{\mathcal{S}}\vdash\theta$ , contrary to the choice of  $\theta$ .

Thus, the same hardness conclusion holds not only for the extension  $\mathcal{S}+Con_{\mathcal{S}}$ , but also for strictly stronger true extensions  $\mathcal{S}+\phi$ . Therefore Pudlák’s Conjecture is not tight as a converse to Theorem 3.4. ■

Thus, there is a gap between Theorem 3.4 and Pudlák’s Conjecture. The former is governed by feasible relative consistency, while the latter already yields hardness for extensions  $\mathcal{S}+\phi$  satisfying  $\mathcal{S}+\phi\vdash Con_{\mathcal{S}}$ , that raise absolute consistency. A central aim of the paper is to isolate principles that close this gap.

### 3.3 Non-Simulation Implies Non-Simulation for Busy Beavers

A theory  $\mathcal{S}$  for which, for every true sentence  $\phi$ , there exists a polynomial bound on  $\mathcal{S}$ -proofs of  $Con_{\mathcal{S}+\phi}(n)$ —would seem to support simulations that cannot be explained solely by information

---

If a chosen formalization yields a true sentence  $\phi_{\text{Jump}}$  such that  $\mathcal{S} + \phi_{\text{Jump}} \vdash Con_{\mathcal{S}}$ , then Pudlák’s Conjecture already implies that  $\mathcal{S}$  does not simulate  $\mathcal{S} + \phi_{\text{Jump}}$ . The stronger **Higher Relative Consistency** assumption yields the same conclusion directly.

provable in  $\mathcal{S}$ , since  $\mathcal{S}$  does not prove all true sentences. We show that this phenomenon is already witnessed by a canonical family of true sentences, namely exact Busy Beaver value statements.

For each  $k$ , let  $t_k=BB(k)$  be the true  $k$ -state Busy Beaver value, and let  $\phi_{BB}(k)$  denote the true sentence asserting the exact value  $BB(k)=t_k$ . Since  $t_k$  is the maximum halting time of any halting  $k$ -state machine on blank input,  $S_2^1$  proves that  $\phi_{BB}(k)$  implies that every halting  $k$ -state machine on blank input halts within  $t_k$  steps. We will use this bounded-halting consequence in the proof of Lemma 3.9.

We begin by proving that, for any fixed true computably axiomatized theory, sufficiently large Busy Beaver axioms already imply its consistency (Aaronson [1, Proposition 3]). The relevant threshold is the size needed to realize the contradiction-search machine for the theory. Once this is in place, any hard true extension of a sound theory  $\mathcal{S}\supseteq S_2^1$  yields eventual hardness throughout the Busy Beaver family.

**Definition 3.8 (Contradiction Search Threshold)** Let  $\mathcal{T}$  be a computably axiomatized theory, and let  $E_{\mathcal{T}}$  be a Turing machine which enumerates  $\mathcal{T}$ -proofs and halts exactly when it finds a proof of contradiction. Define  $k_{\mathcal{T}}^{\text{CS}}$  to be any threshold such that for every  $k\geq k_{\mathcal{T}}^{\text{CS}}$ , there is a  $k$ -state Turing machine computing the same partial function as  $E_{\mathcal{T}}$ .

**Lemma 3.9** *Let  $\mathcal{T}$  be a fixed computably axiomatized true theory extending  $S_2^1$ . Then for every  $k\geq k_{\mathcal{T}}^{\text{CS}}$ , one has  $S_2^1+\phi_{BB}(k)\vdash\text{Con}_{\mathcal{T}}$ .*

**Proof:** Fix  $k\geq k_{\mathcal{T}}^{\text{CS}}$ , and let  $U:=S_2^1+\phi_{BB}(k)$ . By the choice of  $k_{\mathcal{T}}^{\text{CS}}$ , there is a  $k$ -state Turing machine computing the same partial function as  $E_{\mathcal{T}}$ , so we may reason in  $U$  about that  $k$ -state realization of contradiction search for  $\mathcal{T}$ .

Since  $\phi_{BB}(k)$  asserts the exact value  $BB(k)=t_k$ , the theory  $U$  proves that every halting  $k$ -state Turing machine on blank input halts within at most  $t_k$  steps.

Because  $\mathcal{T}$  is true, the actual computation of  $E_{\mathcal{T}}$  does not find a contradiction within the first  $t_k$  steps. This is a fixed finite computation, so by standard bounded-arithmetic formalization of finite computations,  $S_2^1$  proves the corresponding bounded halting fact, and hence so does  $U$ . Therefore  $U\vdash\neg\exists s\leq t_k (E_{\mathcal{T}} \text{ halts in exactly } s \text{ steps})$ .

Combining these two facts,  $U$  proves that the contradiction-search machine for  $\mathcal{T}$  never halts at all, that is,  $U\vdash\neg\exists s (E_{\mathcal{T}} \text{ halts in exactly } s \text{ steps})$ . By the definition of  $E_{\mathcal{T}}$ , this is exactly  $\text{Con}_{\mathcal{T}}$ . Therefore  $S_2^1+\phi_{BB}(k)\vdash\text{Con}_{\mathcal{T}}$ . ■

In particular, if  $\mathcal{T}$  is consistent, then for every  $k\geq k_{\mathcal{T}}^{\text{CS}}$ , the true Busy Beaver sentence  $\phi_{BB}(k)$  is unprovable in  $\mathcal{T}$ , since otherwise  $\mathcal{T}$  would prove  $\text{Con}_{\mathcal{T}}$ .

With this result, we can show that the existence of some hard true extension  $\mathcal{S}+\phi$  is equivalent to eventual hardness throughout the Busy Beaver family:

**Theorem 3.10** *The following are equivalent:*

1. *There exists a true sentence  $\phi$  such that for every constant  $c$ ,  $\mathcal{S}\not\vdash_{\mathcal{T}}^{n^c} \text{Con}_{\mathcal{S}+\phi}(n)$ .*

2. For all sufficiently large  $k$  and every constant  $c$ ,  $\mathcal{S} \not\vdash^{n^c} \text{Con}_{S_2^1 + \phi_{BB}(k)}(n)$ .

**Proof:** (1)→(2) Assume there exists a true sentence  $\phi$  such that for every constant  $c$ ,  $\mathcal{S} \not\vdash^{n^c} \text{Con}_{\mathcal{S} + \phi}(n)$ . Apply Lemma 3.9 with  $\mathcal{T} = \mathcal{S} + \phi$ . Fix sufficiently large  $k$ , and write  $U_k := S_2^1 + \phi_{BB}(k)$ . Then  $U_k \vdash \text{Con}_{\mathcal{S} + \phi}$ .

Fix once and for all a proof  $\pi_k$  of  $\text{Con}_{\mathcal{S} + \phi}$  in  $U_k$ . There is a uniform proof transformation sending any  $(\mathcal{S} + \phi)$ -proof of contradiction to a  $U_k$ -proof of contradiction: given a code of an  $(\mathcal{S} + \phi)$ -proof of  $0=1$  of length at most  $n$ , one appends the fixed derivation  $\pi_k$  of  $\text{Con}_{\mathcal{S} + \phi}$  and the standard verification that the coded object is such a proof. This increases proof length by at most a linear factor. Hence there exist a linear polynomial  $p_k$  and a constant  $d_k$  such that  $\mathcal{S} \vdash^{n^{d_k}} \text{Con}_{U_k}(p_k(n)) \rightarrow \text{Con}_{\mathcal{S} + \phi}(n)$ .

Suppose toward a contradiction that  $\mathcal{S} \vdash^{n^{O(1)}} \text{Con}_{U_k}(n)$  for some sufficiently large  $k$ . Since  $p_k$  is linear, substituting  $p_k(n)$  for  $n$  still gives  $\mathcal{S} \vdash^{n^{O(1)}} \text{Con}_{U_k}(p_k(n))$ . Composing these proofs with this implication yields  $\mathcal{S} \vdash^{n^{O(1)}} \text{Con}_{\mathcal{S} + \phi}(n)$ , contrary to the hypothesis.

Therefore for all sufficiently large  $k$  and every constant  $c$ ,  $\mathcal{S} \not\vdash^{n^c} \text{Con}_{S_2^1 + \phi_{BB}(k)}(n)$ .

(2)→(1) Immediate, since for sufficiently large  $k$  the sentence  $\phi_{BB}(k)$  is true.  $\blacksquare$

Theorem 3.10 yields the following consequence. If a fixed sound theory  $\mathcal{S}$  fails to simulate some true extension  $\mathcal{S} + \phi$ , then for all sufficiently large  $k$  the fixed Busy Beaver family  $(\text{Con}_{S_2^1 + \phi_{BB}(k)}(n))_n$  is already hard for  $\mathcal{S}$ . By itself this gives, for each  $\mathcal{S}$ , only a theory-dependent diagonal family. Under the global hypothesis that every sound theory in the class has some hard true extension, however, one can diagonalize simultaneously over all theory/exponent pairs and obtain a single nonconstructively chosen function  $k(n)$  such that the family  $(\text{Con}_{S_2^1 + \phi_{BB}(k(n))}(n))_n$  is hard for every theory in the class. Via the usual passage from bounded-consistency hardness to propositional hardness, this yields a family of tautologies hard for every proof system arising from such a theory class.

**Theorem 3.11** *Assume that for every sound theory  $\mathcal{S} \supseteq S_2^1$  with polynomial-time decidable axioms there exists a true sentence  $\psi$  such that for every constant  $c$ ,  $\mathcal{S} \not\vdash^{n^c} \text{Con}_{\mathcal{S} + \psi}(n)$ . Then there exists a nonconstructively chosen function  $k: \mathbb{N} \rightarrow \mathbb{N}$  with unbounded range such that for every sound theory  $\mathcal{S} \supseteq S_2^1$  with polynomial-time decidable axioms, the family  $(\text{Con}_{S_2^1 + \phi_{BB}(k(n))}(n))_n$  is hard for  $\mathcal{S}$ . Equivalently, for every such  $\mathcal{S}$  and every constant  $c$ ,  $\mathcal{S} \not\vdash^{n^c} \text{Con}_{S_2^1 + \phi_{BB}(k(n))}(n)$ .*

**Proof:** Let  $(\mathcal{S}_e, d_e)_{e \in \mathbb{N}}$  be an enumeration of all pairs consisting of a sound theory  $\mathcal{S}_e \supseteq S_2^1$  with polynomial-time decidable axioms and a positive integer exponent  $d_e \geq 1$ , with each pair appearing infinitely often.

For each  $e$ , by hypothesis there exists a true sentence  $\psi_e$  such that for every constant  $c$ ,  $\mathcal{S}_e \not\vdash^{n^c} \text{Con}_{\mathcal{S}_e + \psi_e}(n)$ . Therefore, by Theorem 3.10, there exists  $K_e$  such that for every fixed  $m \geq K_e$ ,  $\mathcal{S}_e \not\vdash^{\mathcal{O}(n^{d_e})} \text{Con}_{S_2^1 + \phi_{BB}(m)}(n)$ .

We choose inductively sequences  $m_0 < m_1 < m_2 < \dots$  and  $n_0 < n_1 < n_2 < \dots$ . At stage  $e \geq 1$ , choose  $m_e$  with  $m_e \geq K_e$  and  $m_e > m_{e-1}$ . Since  $\mathcal{S}_e \mid \frac{\mathcal{O}(n^{d_e})}{\mathcal{O}(n^{d_e})} \text{Con}_{S_2^1 + \phi_{BB}(m_e)}(n)$ , there exists  $n_e$  such that every  $\mathcal{S}_e$ -proof of  $\text{Con}_{S_2^1 + \phi_{BB}(m_e)}(n_e)$  has size greater than  $en_e^{d_e}$ . Increase  $n_e$  if necessary so that  $n_e > n_{e-1}$ .

Define  $k(n_e) := m_e$  for each  $e \geq 1$ , and define  $k(n)$  arbitrarily on all other inputs. Since  $(m_e)$  is strictly increasing, the range of  $k$  is unbounded.

Fix a sound theory  $\mathcal{S} \supseteq S_2^1$  with polynomial-time decidable axioms, and suppose toward a contradiction that for some constant  $c$ ,  $\mathcal{S} \mid \frac{n^c}{\mathcal{O}(n^c)} \text{Con}_{S_2^1 + \phi_{BB}(k(n))}(n)$ . Then there exist constants  $A, N$  such that for all  $n \geq N$ , there is an  $\mathcal{S}$ -proof of  $\text{Con}_{S_2^1 + \phi_{BB}(k(n))}(n)$  of size at most  $An^c$ .

Choose  $e$  such that  $\mathcal{S}_e = \mathcal{S}$ ,  $d_e \geq c$ ,  $e \geq A$ , and  $n_e \geq N$ ; this is possible because every theory appears infinitely often paired with arbitrarily large integers. At  $n = n_e$  we have  $k(n_e) = m_e$ , so  $\text{Con}_{S_2^1 + \phi_{BB}(k(n_e))}(n_e) = \text{Con}_{S_2^1 + \phi_{BB}(m_e)}(n_e)$ . By construction every  $\mathcal{S}$ -proof of this sentence has size greater than  $en_e^{d_e}$ , while  $en_e^{d_e} \geq An_e^c$ . This contradicts the assumed upper bound.

Therefore, for every sound theory  $\mathcal{S} \supseteq S_2^1$  with polynomial-time decidable axioms and every constant  $c$ ,  $\mathcal{S} \not\mid \frac{n^c}{\mathcal{O}(n^c)} \text{Con}_{S_2^1 + \phi_{BB}(k(n))}(n)$ . ■

Under Pudlák's Conjecture, one can give a stronger statement indicating just how large  $k$  must be:

**Theorem 3.12** *Assume Pudlák's Conjecture. Then for every  $k \geq k_S^{\text{cs}}$ , for every constant  $c$ ,  $\mathcal{S} \mid \frac{n^c}{\mathcal{O}(n^c)} \text{Con}_{\mathcal{S} + \phi_{BB}(k)}(n)$ .*

**Proof:** Fix  $k \geq k_S^{\text{cs}}$ . Since  $\mathcal{S} + \phi_{BB}(k)$  extends  $S_2^1 + \phi_{BB}(k)$ , Lemma 3.9 implies  $\mathcal{S} + \phi_{BB}(k) \vdash \text{Con}_{\mathcal{S}}$ .

Fix once and for all a proof  $\pi_k$  of  $\text{Con}_{\mathcal{S}}$  in  $\mathcal{S} + \phi_{BB}(k)$ . Since  $k$  is fixed, the size of  $\pi_k$  is a constant independent of  $n$ .

There is a uniform proof transformation sending any  $(\mathcal{S} + \text{Con}_{\mathcal{S}})$ -proof to an  $(\mathcal{S} + \phi_{BB}(k))$ -proof by replacing each use of the extra axiom  $\text{Con}_{\mathcal{S}}$  by the fixed derivation  $\pi_k$ . This increases proof length by at most a constant multiplicative and additive factor. Hence there exist a linear polynomial  $p_k$  and a constant  $d_k$  such that  $\mathcal{S} \mid \frac{n^{d_k}}{\mathcal{O}(n^{d_k})} \text{Con}_{\mathcal{S} + \phi_{BB}(k)}(p_k(n)) \rightarrow \text{Con}_{\mathcal{S} + \text{Con}_{\mathcal{S}}}(n)$ .

Suppose toward a contradiction that  $\mathcal{S} \mid \frac{n^{O(1)}}{\mathcal{O}(n^{O(1)})} \text{Con}_{\mathcal{S} + \phi_{BB}(k)}(n)$ . Since  $p_k$  is linear, substituting  $p_k(n)$  for  $n$  still gives  $\mathcal{S} \mid \frac{n^{O(1)}}{\mathcal{O}(n^{O(1)})} \text{Con}_{\mathcal{S} + \phi_{BB}(k)}(p_k(n))$ . Composing these proofs with the displayed implication yields  $\mathcal{S} \mid \frac{n^{O(1)}}{\mathcal{O}(n^{O(1)})} \text{Con}_{\mathcal{S} + \text{Con}_{\mathcal{S}}}(n)$ , contrary to the hypothesis.

Therefore  $\mathcal{S} \not\mid \frac{n^{O(1)}}{\mathcal{O}(n^{O(1)})} \text{Con}_{\mathcal{S} + \phi_{BB}(k)}(n)$ . Since  $k \geq k_S^{\text{cs}}$  was arbitrary, this holds for every such  $k$ . ■

These results show that Busy Beaver hardness is not merely a convenient encoding of difficult extensions, but reflects a structural barrier to simulation. In each case, non-simulation arises because  $\mathcal{S} + \phi$  carries consistency-strength information whose transfer is not already visible over weak arithmetic.

Informally, any argument showing that a sound theory  $\mathcal{S}$  fails to simulate some true extension  $\mathcal{S}+\phi$  also yields, for that same  $\mathcal{S}$ , some Busy Beaver sentence  $\phi_{BB}(k)$  that  $\mathcal{S}$  cannot prove witnessing the same obstruction. In this sense, the source of non-simulation already appears among canonical incompleteness statements, and reflects a limit on the ability of  $\mathcal{S}$  to exploit true information it cannot itself prove. This is exactly the pattern predicted by **Higher Relative Consistency**: if  $EA \nmid \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ , then  $\text{Con}_{\mathcal{S}+\phi}(n)$  should already be hard for  $\mathcal{S}$ .

## 4 Heuristic Constraints

This section develops several heuristic constraints suggested by the above results. These constraints aim to limit the use of unprovable information in simulations, particularly information drawn from immune sets such as Busy Beavers and Kolmogorov-random strings.

### 4.1 Tightness

A natural heuristic constraint on criteria for non-simulation is to assert that the best known sufficient condition for simulation is already the best possible. Under this heuristic, the sufficient condition from Theorem 3.4 is also necessary, and therefore becomes a characterization of simulation. Equivalently, failure of that condition becomes the predicted criterion for non-simulation. We adopt this only as a working principle: the paper does not prove that Theorem 3.4 is optimal, which would in itself resolve open problems.

### 4.2 Feasible Reflection

A natural structural constraint on simulations is that externally short proofs of bounded consistency should have an internal explanation in a weak base theory. The following definition isolates that requirement.

**Definition 4.1** *Feasible Reflection* holds for  $\mathcal{S}$  if for every true sentence  $\phi \mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$  implies  $EA \nmid \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ .

*Feasible Reflection* asserts that external polynomial-size simulation should already admit a relative-consistency explanation over a weak base theory. In light of Theorem 3.2, when  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}}(n)$ , this means that feasible relative consistency should already imply  $EA \nmid \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ . Informally, the point is not that proving each finite instance should automatically yield a proof of the universal statement, but that efficient proofs of all the bounded consistency statements should not exist unless there is already a weak-base explanation for them. In that sense, *Feasible Reflection* is best viewed not as an induction

principle, but as a structural constraint on efficient proofs that depend on true but unprovable information.

By Theorem 3.2, if  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}}(n)$ , Feasible Reflection states that feasible relative consistency ( $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}}(p(n)) \rightarrow \text{Con}_{\mathcal{S}+\phi}(n)$ ) already yields  $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ .<sup>8</sup>

### 4.3 Simulations Employing Unproven Facts

A natural refinement of the Busy Beaver phenomenon is obtained by considering Kolmogorov-random strings. We fix a universal Turing machine  $U$ . For each string  $x \in \{0, 1\}^*$ , let  $K_U(x)$  denote the plain Kolmogorov complexity of  $x$ , that is, the length of the shortest program  $p$  such that  $U(p) = x$ . Fix once and for all a constant  $c_U$ , and let  $R$  denote the set of strings  $x$  such that  $K_U(x) \geq |x| - c_U$ . Thus  $R$  is the set of Kolmogorov-random strings relative to  $U$ , up to the fixed additive constant  $c_U$ . It is well known that  $R$  is infinite and immune; any theory fails to prove  $x \in R$  with at most finite exceptions, by Chaitin's [3] Incompleteness Theorem (Li and Vitányi [9]).

Kolmogorov-random axioms  $x \in R$  suggest a particularly concrete source of hardness within the relative-consistency viewpoint. The guiding intuition is that efficient simulation should not be able to exploit true random information that the base theory cannot itself recover. This leads to a candidate information-theoretic hardness principle. It should not be viewed as a theorem from the structural discussion above, but rather as a distinguished random-axiom instance of the same general obstruction.

One cannot expect such hardness for every true random axiom. By Theorem 3.4, simulation is possible whenever  $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+(x \in R)}$ . The natural question is therefore whether, in the random-axiom case, the exact obstruction to simulation is failure of access to the axiom already over  $EA + \text{Con}_{\mathcal{S}}$ .

**Conjecture 4.2** (Kolmogorov Hardness) *Whenever  $x \in R$  in the standard model, one has  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}+(x \in R)}(n)$  if and only if  $EA + \text{Con}_{\mathcal{S}} \nmid x \in R$ . Equivalently,  $\mathcal{S} \upharpoonright^{n^{O(1)}} \text{Con}_{\mathcal{S}+(x \in R)}(n)$  if and only if  $EA + \text{Con}_{\mathcal{S}} \vdash x \in R$ .*

**Lemma 4.3** *For each fixed string  $x$ ,  $EA$  proves  $\text{Con}_{\mathcal{S}+(x \in R)} \rightarrow (x \in R)$ . Consequently, if  $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+(x \in R)}$ , then  $EA + \text{Con}_{\mathcal{S}} \vdash x \in R$ .*

---

<sup>8</sup>Freund and Pakhomov [4] give an example, based on PA and slow consistency, showing that external efficient simulation does not by itself force interpretability: PA has polynomial-size proofs of  $\text{Con}_{\text{PA}+\text{Con}^*(\text{PA})}(n)$ , while PA does not interpret  $\text{PA}+\text{Con}^*(\text{PA})$ . This illustrates that additional bridge principles, such as *Feasible Reflection*, are needed to connect simulation with relative-consistency transfer or interpretability.

**Proof:** Fix a string  $x$ . If  $x \notin R$  in the standard model, then choose a specific program  $p$  and a specific number of steps  $t$  such that  $|p| < |x| - c_U$  and  $U(p)$  halts in exactly  $t$  steps with output  $x$ . Since  $x, p, t$  are fixed standard objects,  $EA$  proves that this computation is correct, and hence proves  $\neg(x \in R)$ .

Now  $x \in R$  is an axiom of  $\mathcal{S}+(x \in R)$ . Therefore, from the  $EA$ -proof of  $\neg(x \in R)$ , one can construct a concrete  $\mathcal{S}+(x \in R)$ -proof of contradiction. Because this proof is finite and explicit,  $EA$  proves  $\neg \text{Con}_{\mathcal{S}+(x \in R)}$ .

Thus  $EA$  proves  $\neg(x \in R) \rightarrow \neg \text{Con}_{\mathcal{S}+(x \in R)}$ , and hence  $EA \vdash \text{Con}_{\mathcal{S}+(x \in R)} \rightarrow (x \in R)$ . The consequence follows by composing this implication with  $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+(x \in R)}$ . ■

The point of this formulation is that  $EA + \text{Con}_{\mathcal{S}} \vdash x \in R$  is exactly the kind of information that can make simulation possible. Lemma 4.3 shows that it is also necessary. So the conjecture says that this necessary condition is also sufficient: simulation occurs exactly in those cases where  $EA + \text{Con}_{\mathcal{S}}$  already proves the random axiom.

Moreover, by Chaitin's Incompleteness Theorem applied to the theory  $EA + \text{Con}_{\mathcal{S}}$ , only finitely many true statements of the form  $x \in R$  are provable in that theory. Hence for all sufficiently long Kolmogorov-random strings  $x$ , the condition  $EA + \text{Con}_{\mathcal{S}} \not\vdash x \in R$  automatically holds, and the conjecture therefore predicts  $\mathcal{S} \Big| \frac{n^{O(1)}}{\neq} \text{Con}_{\mathcal{S}+(x \in R)}(n)$ .

## 5 Proposed Characterization of Simulation

Theorem 3.4 shows that relative consistency ( $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ ) implies simulation ( $\mathcal{S} \Big| \frac{n^{O(1)}}{\neq} \text{Con}_{\mathcal{S}+\phi}(n)$ ), for  $\mathcal{S}$  finitely axiomatized and sequential. Taking this theorem to be tight leads to the following structural candidate characterization of simulation. This section develops that proposal and its main consequences. In particular, the case  $\phi = \text{Con}_{\mathcal{S}}$  and the Kolmogorov-randomness axioms arise directly as natural instances of the same general obstruction, while the Busy Beaver reduction shows that arbitrary hard true extensions can be transferred to exact Busy Beaver value statements.

**Assumption 5.1 (Higher Relative Consistency)** *Let  $\mathcal{S} \supseteq \mathcal{S}_2^1$  be a sound theory with polynomial-time decidable axioms. If  $EA \not\vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ , then for every constant  $c > 0$ ,  $\mathcal{S} \Big| \frac{n^c}{\neq} \text{Con}_{\mathcal{S}+\phi}(n)$ . Equivalently, efficient simulation can occur only when the relative-consistency implication is already provable in  $EA$ .*

This assumption is intended as the converse to Theorem 3.4, and thus as a full characterization of simulation in this setting.

### 5.1 Consequences of Higher Relative Consistency

The next theorem summarizes the main properties of **Higher Relative Consistency** and shows that it performs well against the constraints developed earlier. It is tight relative to the

best-known simulation result (Theorem 3.4); it is equivalent to *Feasible Reflection*; and, in the finitely axiomatized sequential case, it implies Pudlák’s Conjecture and *Kolmogorov Hardness*.

**Theorem 5.2** *Let  $\mathcal{S} \supseteq \mathcal{S}_2^1$  be a sound theory with polynomial-time decidable axioms such that  $\mathcal{S} \not\vdash^{n^{O(1)}} \text{Con}_{\mathcal{S}}(n)$ . Then:*

1. *The principles **Higher Relative Consistency** and Feasible Reflection are equivalent.*
2. ***Higher Relative Consistency** implies Pudlák’s Conjecture: for every constant  $c$ ,  $\mathcal{S} \not\vdash^{n^c} \text{Con}_{\mathcal{S}+\text{Con}_{\mathcal{S}}}(n)$ .*

*If, in addition,  $\mathcal{S}$  is finitely axiomatized and sequential, then:*

3. ***Higher Relative Consistency** is tight relative to Theorem 3.4:  $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$  if and only if  $\mathcal{S} \not\vdash^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ .*
4. ***Higher Relative Consistency** implies Kolmogorov Hardness.*

**Proof:** For part (1), first assume **Higher Relative Consistency**. Suppose  $\mathcal{S} \not\vdash^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ . If  $EA \not\vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ , then **Higher Relative Consistency** would imply  $\mathcal{S} \not\vdash^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ , a contradiction. Hence  $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ . Thus *Feasible Reflection* holds.

Conversely, assume *Feasible Reflection*. Suppose  $EA \not\vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ . If  $\mathcal{S} \not\vdash^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ , then by *Feasible Reflection* one would have  $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ , a contradiction. Therefore  $\mathcal{S} \not\vdash^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ . This is exactly **Higher Relative Consistency**.

For part (2), assume **Higher Relative Consistency**. Let  $U := \mathcal{S} + \text{Con}_{\mathcal{S}}$ . Since  $\mathcal{S}$  is sound,  $U$  is consistent and recursively axiomatizable. By Gödel’s Second Incompleteness Theorem,  $U \not\vdash \text{Con}_U$ . But  $\text{Con}_U$  is exactly  $\text{Con}_{\mathcal{S}+\text{Con}_{\mathcal{S}}}$ . Therefore  $EA \not\vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\text{Con}_{\mathcal{S}}}$ : otherwise, because  $U$  extends  $EA$  and contains  $\text{Con}_{\mathcal{S}}$ , the theory  $U$  would prove  $\text{Con}_{\mathcal{S}+\text{Con}_{\mathcal{S}}} = \text{Con}_U$ , contradicting Gödel’s theorem. Hence **Higher Relative Consistency** yields that for every constant  $c > 0$ ,  $\mathcal{S} \not\vdash^{n^c} \text{Con}_{\mathcal{S}+\text{Con}_{\mathcal{S}}}(n)$ .

For part (3), assume **Higher Relative Consistency**. If  $\mathcal{S}$  is finitely axiomatized and sequential and  $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ , then Theorem 3.4 gives  $\mathcal{S} \not\vdash^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ .

Conversely, suppose  $\mathcal{S} \not\vdash^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ . If  $EA \not\vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ , then **Higher Relative Consistency** would imply  $\mathcal{S} \not\vdash^{n^{O(1)}} \text{Con}_{\mathcal{S}+\phi}(n)$ , a contradiction. Therefore  $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+\phi}$ .

For part (4), assume **Higher Relative Consistency**, and suppose  $\mathcal{S}$  is finitely axiomatized and sequential. Let  $x \in R$  be true.

First suppose  $\mathcal{S} \not\vdash^{n^{O(1)}} \text{Con}_{\mathcal{S}+(x \in R)}(n)$ . By part (1), this implies  $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow \text{Con}_{\mathcal{S}+(x \in R)}$ . By Lemma 4.3,  $EA$  proves  $\text{Con}_{\mathcal{S}+(x \in R)} \rightarrow (x \in R)$ . Hence  $EA + \text{Con}_{\mathcal{S}} \vdash x \in R$ .

Conversely, suppose  $EA + \text{Con}_{\mathcal{S}} \vdash x \in R$ . By deduction,  $EA \vdash \text{Con}_{\mathcal{S}} \rightarrow (x \in R)$ . Since  $EA$  can formalize the trivial proof transformation replacing uses of the extra axiom  $(x \in R)$  by a proof of

$(x \in R)$  obtained from  $Con_{\mathcal{S}}$ , it follows that  $EA \vdash Con_{\mathcal{S}} \rightarrow Con_{\mathcal{S}+(x \in R)}$ . Since  $\mathcal{S}$  is finitely axiomatized and sequential, part (3) (equivalently, Theorem 3.4) now yields  $\mathcal{S} \not\vdash^{n^{O(1)}} Con_{\mathcal{S}+(x \in R)}(n)$ .

Therefore  $\mathcal{S} \not\vdash^{n^{O(1)}} Con_{\mathcal{S}+(x \in R)}(n)$  if and only if  $EA + Con_{\mathcal{S}} \vdash x \in R$ . Equivalently,  $\mathcal{S} \not\vdash^{n^{O(1)}} Con_{\mathcal{S}+(x \in R)}(n)$  if and only if  $EA + Con_{\mathcal{S}} \not\vdash x \in R$ . This is exactly *Kolmogorov Hardness*. ■

## 6 Further Conjectural Extensions

The core paper ends with **Higher Relative Consistency** and Conjecture 4.2. The present section presents strictly stronger conjectures motivated by that picture, where hardness emerges for small  $n$  and not only asymptotically. These are not consequences of the previous section and are not needed for the main characterization program. Their role is to indicate what additional phenomena would follow if the random-axiom obstruction were strengthened from polynomial simulation hardness to subexponential hardness at the first syntactically meaningful proof scale.

Two additional choices are needed. First, one must specify the target lower bound on proof length; here we take subexponential hardness against proofs of size  $2^{(1-\epsilon)n}$ . Second, one must specify the least scale at which the relevant bounded-consistency statement becomes genuinely nontrivial. Conceptually, there are two constraints. One is semantic: the ambient string length  $n$  should lie beyond the point at which the base theory no longer already proves the randomness statements under consideration. The other is syntactic: the bounded-consistency parameter should be large enough that a proof of length at most  $n$  can even mention the added axiom. For very small proof bounds the statement may be vacuous, and even after vacuity disappears the bound may still be too small to mention the added axiom. In addition, we include one harmless baseline term recording the least length of a valid proof string in the fixed proof calculus. We therefore package all of these into a single threshold.

To make this explicit, fix a base theory  $\mathcal{S}$ . Let  $N_R^{\mathcal{S}}$  be the maximum of the following three quantities: the least length of a valid  $\mathcal{S}$ -proof string in the fixed proof calculus; a semantic threshold beyond which  $\mathcal{S}$  proves no sentence of the form  $x \in R$ ; and a syntactic threshold beyond which a proof of length at most  $n$  can already mention the added axiom  $x \in R$ .

The easy direction in the finite-scale principles below is routine.

**Lemma 6.1** *Assume  $\mathcal{S} \not\vdash^{n^{O(1)}} Con_{\mathcal{S}}(n)$ . If  $\mathcal{S} \vdash \phi$ , then  $\mathcal{S} \not\vdash^{n^{O(1)}} Con_{\mathcal{S}+\phi}(n)$ .*

**Proof:** Since  $\mathcal{S} \vdash \phi$ , any proof in  $\mathcal{S}+\phi$  can be converted into a proof in  $\mathcal{S}$  by replacing each use of the extra axiom  $\phi$  by a fixed  $\mathcal{S}$ -proof of  $\phi$ . Weak arithmetic formalizes this proof transformation with only polynomial overhead. Hence polynomial-size  $\mathcal{S}$ -proofs of  $Con_{\mathcal{S}}(n)$  yield polynomial-size  $\mathcal{S}$ -proofs of  $Con_{\mathcal{S}+\phi}(n)$ . ■

By the choice of  $N_R^S$ , Lemma 6.1 gives the easy direction in the conjecture below: whenever  $\mathcal{S} \vdash x \in R$  and  $n > N_R^S$ , the bounded-consistency statement  $\text{Con}_{\mathcal{S}+x \in R}(n)$  has polynomial-size, and hence also subexponential-size,  $\mathcal{S}$ -proofs.

**Conjecture 6.2** (*SETH-K-Finite*) For every  $\epsilon > 0$ , every  $n > N_R^S$ , and every string  $x \in \{0, 1\}^n$ , one has  $\mathcal{S} \mid_{\frac{2^{(1-\epsilon)n}}{7}} \text{Con}_{\mathcal{S}+x \in R}(n)$  if and only if  $\mathcal{S} \vdash x \in R$ .

## 6.1 No Mutual Help

*Kolmogorov Hardness* alone does not imply a no-mutual-help statement formulated in terms of provability in the extended theory  $\mathcal{T} := \mathcal{S} + x \in R$ . Applied to the base theory  $\mathcal{T}$ , it yields at most a criterion in terms of whether  $EA + \text{Con}_{\mathcal{T}} \vdash y \in R$ , not in terms of whether  $\mathcal{T} \vdash y \in R$ . **SETH-K-Finite** is therefore a strictly stronger finite-scale strengthening, not a consequence of the previous section. This says that, above the relevant threshold for  $R$ , the only subexponentially easy bounded-consistency instances arising from exact randomness axioms are those whose axioms were already theorems of the base theory.

Now let  $x \in R$  be true and put  $\mathcal{T} := \mathcal{S} + x \in R$ . Let  $N_R^{\mathcal{T}}$  be the threshold defined for the base theory  $\mathcal{T}$  and the predicate  $R$  in the same way that  $N_R^S$  was defined for  $\mathcal{S}$ .

**Theorem 6.3** Assume *SETH-K-Finite*. Let  $x \in R$  be true, let  $\mathcal{T} := \mathcal{S} + x \in R$ , and suppose  $\mathcal{T}$  is sound and belongs to the same class of theories as  $\mathcal{S}$ . Then for every  $\epsilon > 0$ , every  $n > N_R^{\mathcal{T}}$ , and every string  $y \in \{0, 1\}^n$ , one has  $\mathcal{T} \mid_{\frac{2^{(1-\epsilon)n}}{7}} \text{Con}_{\mathcal{T}+y \in R}(n)$  if and only if  $\mathcal{T} \vdash y \in R$ .

**Proof:** Apply **SETH-K-Finite** with base theory  $\mathcal{T}$ . ■

**Corollary 6.4** Assume *SETH-K-Finite*. Let  $x \in R$  be true, let  $\mathcal{T} := \mathcal{S} + x \in R$ , and suppose  $\mathcal{T}$  is sound and belongs to the same class of theories as  $\mathcal{S}$ . Then for every  $\epsilon > 0$ , every  $n > N_R^{\mathcal{T}}$ , and every string  $y \in \{0, 1\}^n$ , if  $\mathcal{T} \not\vdash y \in R$ , then  $\mathcal{T} \mid_{\frac{2^{(1-\epsilon)n}}{7}} \text{Con}_{\mathcal{T}+y \in R}(n)$ .

**Proof:** This is the contrapositive direction of Theorem 6.3. ■

## 6.2 Density of Hard Sentences

For density arguments, the set  $R$  is not the most convenient one, because it need not occupy a large fraction of  $\{0, 1\}^n$ . For that purpose it is better to pass to the logarithmic-deficiency set  $R^{\log}$ , defined by  $x \in R^{\log}$  if and only if  $K_U(x) \geq |x| - d \log |x|$ . Below, we will refer to **SETH-K-Finite** with  $R$  replaced by  $R^{\log}$ ; this yields a finite-scale strengthening on a set of strings whose density tends to 1.

**Theorem 6.5** Assume *SETH-K-Finite* (with  $R$  replaced with  $R^{\log}$ ). Then for every  $\epsilon > 0$  and every  $n > N_{R^{\log}}^S$ , at least  $(1 - O(n^{-d}))2^n$  strings  $x \in \{0, 1\}^n$  satisfy  $\mathcal{S} \mid_{\frac{2^{(1-\epsilon)n}}{7}} \text{Con}_{\mathcal{S}+x \in R^{\log}}(n)$ .

**Proof:** By definition of  $R^{\log}$ , if  $x \notin R^{\log}$  then  $K_U(x) < |x| - d \log |x|$ . Hence the number of strings of length  $n$  outside  $R^{\log}$  is at most  $2^{n-d \log n+1} = O(2^n/n^d)$ . So at least  $(1-O(n^{-d}))2^n$  strings of length  $n$  lie in  $R^{\log}$ .

Now let  $x \in R^{\log} \cap \{0, 1\}^n$ . Since  $n > N_{R^{\log}}^{\mathcal{S}}$ , the defining property of  $N_{R^{\log}}^{\mathcal{S}}$  gives  $\mathcal{S} \not\models x \in R^{\log}$ . By **SETH-K-Finite** (for  $R^{\log}$ ), it follows that  $\mathcal{S} \not\models \frac{2^{(1-\epsilon)n}}{n} \text{Con}_{\mathcal{S}+x \in R^{\log}}(n)$ . ■

Thus, for every  $n > N_{R^{\log}}^{\mathcal{S}}$ , hard bounded-consistency instances occupy a  $(1-O(n^{-d}))$  fraction of  $\{0, 1\}^n$ . If these bounded-consistency statements are translated into tautologies in the usual way, then the resulting tautologies are hard on a density- $1-O(n^{-d})$  set at each such length.

## 7 Remarks on Provability

A complementary question is whether **Higher Relative Consistency** might itself be unprovable. This section is exploratory: it does not establish an independence result for **Higher Relative Consistency**, but records reasons one might expect that natural uniform formulations of the principle will be difficult to prove inside weak first-order theories. The first reason is model-theoretic: the intended meaning of **Higher Relative Consistency** is external and specific to the standard model, whereas any proof in weak arithmetic must proceed through the internal arithmetized proof predicate. The second is complexity-theoretic: a sufficiently uniform proof of **Higher Relative Consistency** would already yield the major lower bounds that motivate the paper.

The model-theoretic issue is basic. Let  $\text{Prf}_{\mathcal{S}}(p, q)$  be the usual formula expressing that  $p$  codes an  $\mathcal{S}$ -proof of the sentence with Gödel number  $q$ . If  $M \models \mathcal{S}$  is nonstandard, then an element  $p \in M$  may satisfy  $M \models \text{Prf}_{\mathcal{S}}(p, \ulcorner \varphi \urcorner)$  even though  $p$  is nonstandard and does not correspond to any genuine finite proof in the standard model. From the internal point of view of  $M$ , however, such a  $p$  is simply a proof of  $\varphi$ . Thus nonstandard models can contain spurious witnesses to provability that have no external proof-theoretic meaning.

This matters because the assertion that  $\mathcal{S} \not\models \frac{n^c}{n} \text{Con}_{\mathcal{S}+\phi}(n)$  is intended externally. It means that for each standard  $n$  there is a genuine finite  $\mathcal{S}$ -proof of the standard sentence  $\text{Con}_{\mathcal{S}+\phi}(n)$  whose length is bounded by a fixed polynomial in  $n$ . Any internal first-order formalization, however, necessarily replaces this by quantification over all elements of a model, including nonstandard proof codes and nonstandard bounds. Thus the internalized sentence can hold in a nonstandard model for reasons that have nothing to do with genuine external simulation.

Accordingly, **Higher Relative Consistency** is not merely a statement about the absence of internal witnesses to short proofs. It is a statement that internal proof data correctly tracks external simulation data in the standard model. But that is precisely what first-order arithmetic cannot in general enforce from within: it has no way to isolate the standard proof codes from nonstandard ones, nor the standard proof-length bounds from nonstandard ones. In this sense, **Higher Relative Consistency** behaves less like an ordinary reflection principle

and more like a standardness principle saying that only genuine external polynomial-size proofs should count as simulation data.

The same issue appears in the random-axiom setting. Even when  $\mathcal{S} \not\vdash x \in R$ , a nonstandard model may still internally treat the bounded-consistency statements for  $\mathcal{S} + (x \in R)$  as having short  $\mathcal{S}$ -proofs, because nonstandard proof codes can witness the corresponding arithmetized proof predicates. Thus an internalized version of the simulation claim may hold in such a model even when no genuine polynomial simulation exists externally. This is exactly why principles such as **Higher Relative Consistency** and *Kolmogorov Hardness* are best understood as external constraints on when apparent simulation data should be taken seriously.

There is also a second reason to doubt provability. A sufficiently uniform proof of **Higher Relative Consistency** would not be a modest metamathematical tidying-up result. Combined with the unconditional simulation theorem and the Busy Beaver reduction, it would already yield sweeping lower bounds, including canonical Busy Beaver witnesses to non-simulation and, in the global form, the nonexistence of an optimal theory or proof system. Thus any proof of **Higher Relative Consistency** strong enough to support the paper's main applications would already settle central open problems.

This point is especially clear in the uniform setting. For a fixed pair  $(\mathcal{S}, \phi)$ , an instance of **Higher Relative Consistency** can be arithmetized as an ordinary first-order sentence about proof predicates and proof lengths. Thus, for a specific theory such as *ZFC*, it is meaningful to ask whether that particular instance is provable, refutable, or independent. However, the discussion above does not establish such an independence result for any fixed instance. Rather, it explains why the most natural and useful forms of **Higher Relative Consistency** are not merely ordinary internal first-order assertions.

The difficulty increases once one passes from fixed instances to uniform forms. A theorem ranging over all true  $\phi$ , even for one fixed  $\mathcal{S}$ , must still relate internal proof-predicate assertions to genuine standard-model simulation behavior. That is exactly where nonstandard models interfere. What is difficult is not writing down the arithmetization, but justifying the passage from internal proof data to external polynomial simulation claims. The difficulty becomes sharper still in the fully global form ranging over all sound  $\mathcal{S}$  and all true  $\phi$ , since both soundness and truth are themselves external semantic conditions.

This also clarifies the role of stronger foundational frameworks. Merely moving from first-order to second-order syntax does not by itself remove the problem, so long as one remains inside a formal theory with its own internal proof predicate and associated nonstandard proof theory. The real issue is not the order of the language, but access to standard truth. For exact Busy Beaver sentences  $\phi_{BB}(k)$ , the obstruction is already visible: a uniform treatment of all such sentences requires more than the resources needed merely to formalize bounded consistency. The upper-bound side of an exact Busy Beaver statement is naturally  $\Pi_1$ , while the lower-bound side is witnessed by an explicit finite halting computation. Thus a metatheory with a truth predicate for all true  $\Pi_1$  sentences, together with ordinary arithmetic verification of finite computations, is a natural setting for uniform reasoning about the Busy Beaver family.

For **Higher Relative Consistency** in full generality, ranging over arbitrary true  $\phi$ , one would want a broader truth or satisfaction framework, or else an explicitly external axiom schema.

Accordingly, the right conclusion is not that the preceding argument proves independence of **Higher Relative Consistency** from *ZFC* or from any other specific theory. Rather, it shows that the global form of **Higher Relative Consistency** is most naturally understood as an external structural principle whose faithful formulation already points beyond ordinary first-order arithmetic. If **Higher Relative Consistency** is true, its role may therefore be closer to that of an external axiom governing when apparent simulation data reflects genuine proof-theoretic structure.

## 8 Conclusion

This paper studies the Characterization Problem for simulation between arithmetic theories and advances a specific conjectural answer to its hardness side. It does not solve the problem. Rather, its unconditional results isolate constraints that any successful characterization should satisfy, while its conjectural part proposes one way those constraints might fit together.

The paper's two unconditional contributions are structural. First, feasible relative consistency suffices for simulation, giving a general upper-bound mechanism that includes the case in which  $EA \vdash Con_S \rightarrow Con_{S+\phi}$ . Second, any hard true extension can be replaced by one of Busy Beaver form, showing that canonical incompleteness phenomena already suffice to witness non-simulation.

Against this backdrop, the paper proposes **Higher Relative Consistency** as its central structural criterion: if  $EA \not\vdash Con_S \rightarrow Con_{S+\phi}$ , then for every constant  $c > 0$ ,  $\mathcal{S} \not\vdash^{nc} Con_{S+\phi}(n)$ . This is the paper's main candidate characterization of simulation in terms of relative-consistency transfer in a weak base theory.

Within this framework, the paper also formulates *Kolmogorov Hardness* as a distinguished random-axiom instance of the same general obstruction. Its guiding idea is that canonical random axioms should instantiate the hardness predicted by failure of weak-base relative-consistency transfer. In this sense, *Kolmogorov Hardness* is not a competing criterion, but a particularly natural specialization of **Higher Relative Consistency**.

Proving **Higher Relative Consistency** would amount to a major breakthrough. As the missing converse to the strongest known upper-bound mechanism, it would convert relative-consistency transfer in a weak base theory from a sufficient condition into a genuine characterization of simulation. Whether that proposal is correct remains open, but the present paper substantially narrows the space of plausible alternatives: any successful characterization must explain sufficient conditions for simulation and why canonical incompleteness phenomena, already at the level of Busy Beaver truths and random axioms, suffice to witness its failure.

## References

- [1] Scott Aaronson, *The busy beaver frontier*, SIGACT News **51** (2020), no. 3, 32–54.
- [2] Samuel R. Buss, *Bounded arithmetic*, Lecture notes, Bibliopolis, 1986.
- [3] Gregory J. Chaitin, *Information-theoretic limitations of formal systems*, JACM **21** (1974), no. 3, 403–424.
- [4] Anton Freund and Fedor Pakhomov, *Short proofs for slow consistency*, Notre Dame Journal of Formal Logic **61** (2020), no. 1, 31–49.
- [5] Petr Hájek and Pavel Pudlák, *Metamathematics of first-order arithmetic*, Perspectives in Logic, Cambridge University Press, 1998.
- [6] Erfan Khaniki, *Jump operators, interactive proofs and proof complexity generators*, 2024 IEEE 65th Annual Symposium on Foundations of Computer Science (FOCS), 2024, pp. 573–593.
- [7] Jan Krajíček, *Proof complexity*, Cambridge University Press, New York, NY, 2019.
- [8] Jan Krajíček and Pavel Pudlák, *Propositional proof systems, the consistency of first order theories and the complexity of computations*, J. Symb. Log. **54** (1989), 1063–79.
- [9] Ming Li and Paul M. B. Vitányi, *An introduction to Kolmogorov complexity and its applications*, Texts in Computer Science, Springer, 2008.
- [10] Pavel Pudlák, *On the length of proofs of finitistic consistency statements in first order theories*, Studies in Logic and the Foundations of Mathematics, vol. 120, Elsevier, 1986, pp. 165–196.
- [11] ———, *Logical foundations of mathematics and computational complexity: A gentle introduction*, Springer, 2013.
- [12] ———, *Incompleteness in the finite domain*, Bull. Symb. Log. **23** (2017), no. 4, 405–441.
- [13] Tibor Rado, *On non-computable functions*, The Bell System Technical Journal **41** (1962), no. 3, 877–884.
- [14] Albert Visser, *The interpretation existence lemma*, Feferman on Foundations: Logic, Mathematics, Philosophy (Gerhard Jäger and Wilfried Sieg, eds.), Springer International Publishing, Cham, 2017, pp. 101–144.