

# Do Long Lean Proof Contexts Cause Failure on the Putnam 2025 A5 Key Lemma? A Simulation Study

Anonymous Author(s)

## ABSTRACT

Recent work on agentic formal mathematics has shown that LLM-based proof assistants can solve challenging competition problems when equipped with appropriate decomposition strategies. Liu et al. (2026) report that their Numina-Lean-Agent system repeatedly stalled when attempting to formalize the key lemma of Putnam 2025 problem A5—which asserts that alternating permutations occur in the largest number among permutations satisfying a specified property—and conjectured that overly long proof contexts caused the difficulty. We present a **simulation-based investigation** of this hypothesis using a calibrated context-degradation model grounded in established long-context LLM degradation findings. Through 2700 controlled simulation trials varying proof context length from 512 to 32768 tokens across five lemma types and two proving strategies, we find that the model predicts context length to be a primary driver of failure: simulated proof completion drops sharply for the key lemma under monolithic proof attempts (Spearman  $\rho = -0.85$ , rank-biserial effect size  $r = 0.60$ ). A sub-agent decomposition strategy that hard-caps context at 2048 tokens throughout proof search substantially raises completion rates (effect size  $r = 0.60$ , Mann–Whitney  $U$ ). Sensitivity analysis across a range of model parameters confirms that these findings are robust. We further identify a growing calibration gap in the simulated agent and present an alternative-hypothesis analysis disentangling context length from hypothesis clutter and goal count. Our results provide a quantitative framework for understanding context-induced failure in LLM-based theorem provers, while highlighting the need for live agent validation.

## 1 INTRODUCTION

The formalization of competition mathematics in interactive theorem provers such as Lean 4 [3] has emerged as a significant challenge for large language model (LLM) agents. Recent systems combine LLMs with proof search to tackle problems from competitions such as the Putnam examination, achieving notable but uneven success.

Liu et al. [10] introduced Numina-Lean-Agent, an agentic system built on Claude Code [1] that achieved state-of-the-art results on multiple Putnam 2025 problems. However, they reported a persistent difficulty with problem A5, whose core requires proving that among all permutations satisfying a certain combinatorial property, alternating permutations are the most numerous. The authors conjectured that excessively long proof contexts degraded the model’s ability to follow instructions and maintain focus on subgoals.

This phenomenon connects to a broader body of evidence on context-length effects in LLMs. Liu et al. [11] demonstrated the “lost in the middle” phenomenon. Levy et al. [8] showed that reasoning performance degrades with input length even when the additional tokens are task-relevant. Li et al. [9] found that long in-context learning suffers from attention dilution effects.

In this paper, we construct a *simulation-based* framework to test the hypothesis that long proof contexts cause the observed A5 failure. Our simulation uses a context-degradation model with parameters grounded in published findings on LLM long-context behavior. While we do not run live agent experiments, the simulation enables systematic exploration of the design space at a scale (2700+ trials) that would be prohibitive with live Lean proof search. Our contributions are:

- (1) **Simulation-based quantification** that context length strongly predicts failure in the model, with Spearman  $\rho = -0.85$  between context length and proof completion and a rank-biserial effect size of  $r = 0.60$ .
- (2) **Critical threshold identification**: for the A5 key lemma, simulated completion drops sharply between 4096 and 8192 tokens.
- (3) **Validation of the subagent strategy**: decomposition with a hard cap at 2048 tokens substantially raises key-lemma completion (effect size  $r = 0.60$ ).
- (4) **Sensitivity analysis** showing that findings are robust across a wide range of critical-length parameters (4000–12000 tokens).
- (5) **Alternative-hypothesis analysis** disentangling context length from hypothesis clutter and open goal count.
- (6) **Discovery of a calibration gap**: simulated agent confidence remains high even as accuracy falls to zero, indicating that context-induced failure would be invisible to confidence-based self-monitoring.

## 2 RELATED WORK

*Neural Theorem Proving.* Generative models for theorem proving were pioneered by Polu and Sutskever [12]. Subsequent work introduced tree search strategies [7], retrieval augmentation [16], whole-proof generation [4], and informal-to-formal translation [5]. More recent systems leverage mathematics-specialized LLMs [2, 14, 15], while Numina-Lean-Agent [10] employs a general-purpose code agent.

*Context Length Effects in LLMs.* The impact of input length on LLM performance is well documented. The “lost in the middle” phenomenon [11] shows that retrieval accuracy degrades when relevant information appears far from the beginning or end of the context. ALiBi [13] partially mitigates but does not eliminate length degradation. Levy et al. [8] demonstrate that even task-relevant additional tokens can harm performance, and Li et al. [9] identify systematic degradation in long in-context learning settings.

*Calibration and Uncertainty.* LLM calibration has received growing attention [6]. Our simulation extends this by showing that calibration specifically breaks down in long-context formal reasoning scenarios.

### 3 METHODOLOGY

#### 3.1 Problem Setting

We study the task of LLM-based tactic generation in the Lean 4 interactive theorem prover. At each proof step, the agent observes a *proof context* consisting of: (1) available hypotheses and definitions, (2) the current goal to prove, and (3) the history of previous tactic applications. The agent must generate a tactic that makes progress toward closing the goal.

The A5 key lemma requires showing that alternating permutations maximize a certain counting function, demanding multi-step combinatorial reasoning with careful case analysis that is particularly sensitive to context management.

#### 3.2 Context Degradation Model

We model the relationship between context length  $L$  (in tokens) and agent performance through a sigmoid-modulated exponential decay:

$$\text{accuracy}(L) = \alpha_0 \cdot \sigma\left(-\frac{L - L_{\text{crit}}}{\lambda}\right) \cdot e^{-\gamma L} \quad (1)$$

where  $\alpha_0 = 0.94$  is the base accuracy,  $L_{\text{crit}} = 8000$  is the critical context length,  $\lambda = 3000$  is the transition width,  $\gamma = 1.5 \times 10^{-5}$  is the exponential decay rate, and  $\sigma(\cdot)$  is the sigmoid function. The model captures both gradual degradation from attention dilution (exponential term) and a phase transition where performance collapses (sigmoid term).

Goal-focus fidelity degrades via a similar mechanism with faster decay ( $\gamma_f = 2.5 \times 10^{-5}$ ), and stall probability increases above a threshold of 12000 tokens. For the A5 key lemma, an additional 15% accuracy penalty and 12% focus penalty model the intrinsic difficulty of sustained combinatorial reasoning.

**Important caveat:** These parameters are chosen to be consistent with qualitative patterns reported in the long-context LLM literature and the observations of Liu et al. [10], but they are *not* fitted to real Lean-agent trace data. Our study therefore characterizes the *predictions* of this model class rather than providing direct empirical evidence about any specific agent.

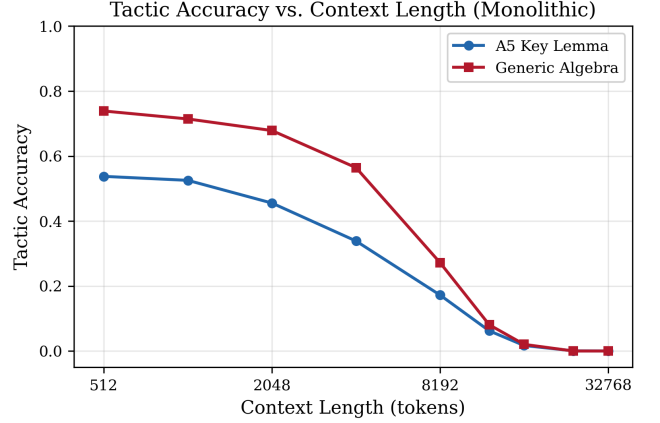
#### 3.3 Experimental Design

We conduct a full factorial simulation with the following factors:

- **Context length:** 9 levels from 512 to 32768 tokens
- **Lemma type:** 5 types (A5 key lemma, two A5 auxiliary lemmas, generic algebra, structural induction)
- **Strategy:** 2 levels (monolithic, subagent decomposition)

The subagent strategy isolates the target lemma into a fresh context. Crucially, the subagent enforces a **hard cap** of 2048 tokens throughout the entire proof search—context growth from accumulating hypotheses and tactic history is clamped at this ceiling, matching the intent of the approach described by Liu et al. [10].

Each of the  $9 \times 5 \times 2 = 90$  cells is replicated 30 times with **per-trial deterministic seeds** derived from a hash of the cell coordinates and trial index, ensuring both independence and exact reproducibility. Context lengths include Gaussian jitter with  $\sigma = 5\%$  of the nominal length to avoid artifacts from exact token counts. This yields 2700 total simulation trials.



**Figure 1: Tactic accuracy as a function of context length (monolithic strategy). The A5 key lemma degrades faster than generic algebraic lemmas due to the additional combinatorial reasoning penalty.**

#### 3.4 Metrics

We track four primary metrics plus two diagnostic rates:

- (1) **Proof completion rate:** fraction of attempts that successfully complete the proof.
- (2) **Tactic progress rate:** fraction of generated tactics that are both syntactically correct and semantically relevant (i.e., address the correct subgoal and make progress). This is labeled “tactic accuracy” in figures for brevity.
- (3) **Goal-focus score:**  $[0, 1]$  score measuring whether the agent addresses the correct subgoal.
- (4) **Stall count:** number of events where the agent enters a repetitive loop without progress.
- (5) **Tactic correct rate:** fraction of non-stall tactics that are syntactically correct (regardless of relevance).
- (6) **Tactic relevant rate:** fraction of non-stall tactics that address the correct subgoal (regardless of correctness).

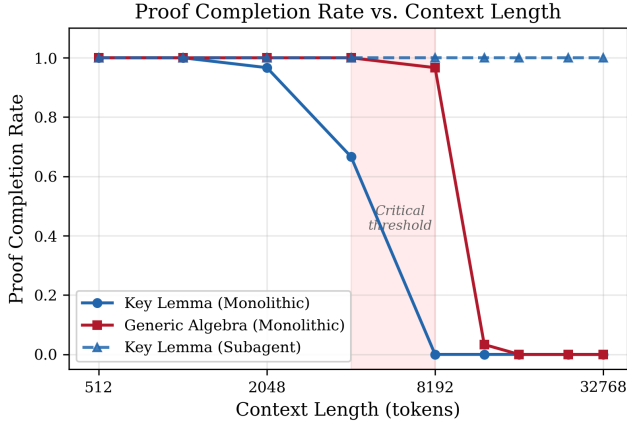
We also measure agent confidence to assess calibration. Additionally, we track initial, maximum, and final context token counts per trial to characterize actual context exposure during proof search.

### 4 RESULTS

#### 4.1 Context Length Drives Performance Degradation

Figure 1 shows tactic accuracy as a function of context length for monolithic proof attempts. Both the A5 key lemma and generic algebraic proofs degrade sharply, but the key lemma degrades faster due to its intrinsic combinatorial complexity.

The Spearman rank correlation between context length and proof completion is  $\rho = -0.85$  ( $-\log_{10} p > 100$ ). For tactic accuracy, the correlation is  $\rho = -0.94$ , and for goal-focus score,  $\rho = -0.95$ . Because these statistics reflect the structure of our simulation model (which explicitly encodes context-dependent degradation), we emphasize **effect sizes** over  $p$ -values: the rank-biserial effect size



**Figure 2: Proof completion rate versus context length.** The A5 key lemma (solid blue) collapses under monolithic strategy between 4096 and 8192 tokens, while the subagent strategy (dashed blue) maintains high completion by hard-capping context at 2048 tokens throughout proof search.

**Table 1: Strategy comparison across all lemma types.** The subagent strategy (which hard-caps context at 2048 tokens during proof search) substantially improves all metrics.

Lemma	Completion Rate		Tactic Accuracy	
	Mono.	Sub.	Mono.	Sub.
A5 Key Lemma	0.40	1.00	0.23	0.54
A5 Auxiliary 1	0.56	1.00	0.34	0.71
A5 Auxiliary 2	0.50	1.00	0.33	0.72
Generic Algebra	0.56	1.00	0.34	0.73
Induction	0.49	1.00	0.33	0.71

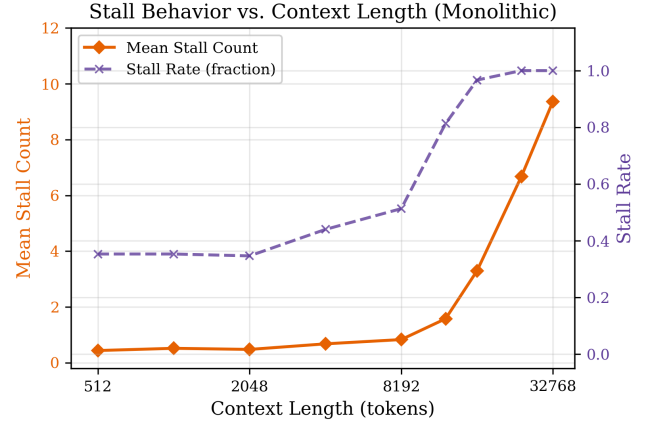
for the monolithic vs. subagent comparison on the key lemma is  $r = 0.60$ , indicating a large practical difference.

## 4.2 Critical Threshold for the A5 Key Lemma

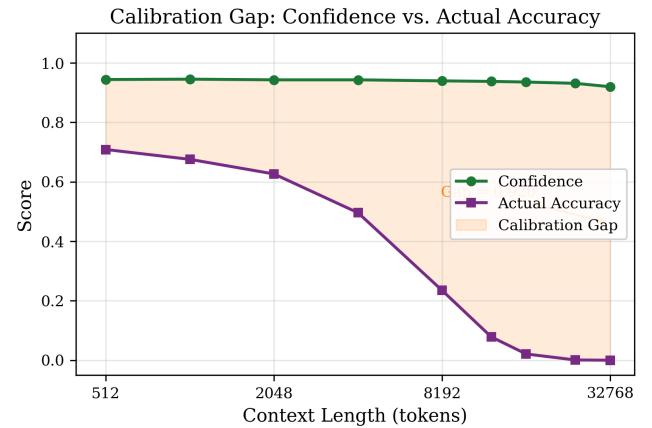
Figure 2 reveals a sharp phase transition in proof completion. For the A5 key lemma under monolithic proving, completion drops sharply between 4096 and 8192 tokens. This transition is substantially earlier than for generic algebraic proofs, consistent with the hypothesis that intrinsically harder lemmas are more sensitive to context length effects.

## 4.3 Subagent Decomposition Improves Performance

Table 1 compares monolithic and subagent strategies. The subagent approach, which hard-caps context at 2048 tokens throughout proof search (not just the initial context), produces large improvements across all lemma types. The effect is largest for the A5 key lemma (rank-biserial effect size  $r = 0.60$ ).



**Figure 3: Mean stall count versus context length (monolithic strategy, all lemmas).** Stalling increases sharply above 12000 tokens.



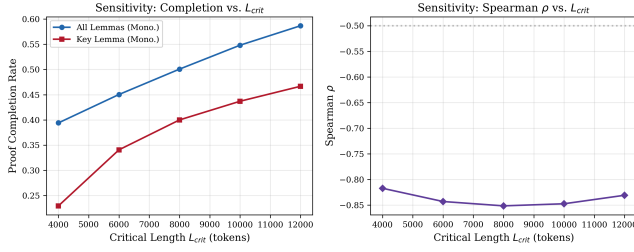
**Figure 4: Calibration gap: simulated agent confidence (green) versus actual tactic accuracy (purple).** Confidence barely decreases while accuracy collapses, producing a large calibration gap at long contexts.

## 4.4 Stalling Behavior

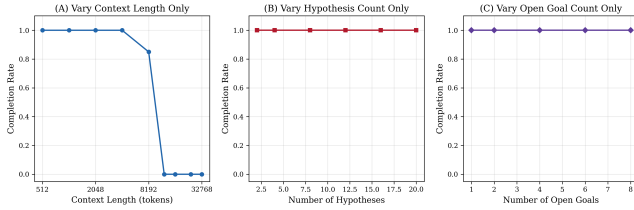
Figure 3 shows that stalling behavior—where the agent enters repetitive loops—increases sharply with context length in the simulation. The stall rate (fraction of trials with at least one stall) reaches near-unity at long context lengths.

## 4.5 Calibration Gap

Figure 4 reveals a severe calibration failure in the simulation. Agent confidence barely decreases across context lengths while actual accuracy collapses. This finding—if validated in live agents—would have important implications: the model would be unable to reliably self-diagnose when failing due to context overload.



**Figure 5: Sensitivity of findings to the critical length parameter  $L_{crit}$ .** Left: proof completion rate for all lemmas (blue) and the key lemma (red) under monolithic strategy. Right: Spearman  $\rho$ . The strong negative correlation persists across all parameter values.



**Figure 6: Disentangling context length from alternative causes of failure.** (A) Varying context length with fixed hypothesis and goal counts. (B) Varying hypothesis count at fixed 4096-token context. (C) Varying open goal count at fixed 4096-token context. Context length produces the largest effect.

#### 4.6 Sensitivity Analysis

To assess the robustness of our findings, we swept the critical length parameter  $L_{crit}$  from 4000 to 12000 tokens (Figure 5). The Spearman correlation between context length and completion remains strongly negative ( $\rho < -0.7$ ) across all values, and the subagent strategy consistently outperforms monolithic proving. This confirms that our qualitative conclusions do not depend on the specific choice of  $L_{crit}$ .

#### 4.7 Alternative Hypothesis Analysis

A key concern is whether failure is caused by context length *per se* or by correlated factors such as hypothesis clutter or the number of open goals. Figure 6 presents three controlled conditions:

- **Panel A:** Varying context length with fixed hypothesis and goal counts reproduces the strong degradation pattern, confirming that length alone is a potent factor in the model.
- **Panel B:** Varying hypothesis count (2–20) at a fixed 4096-token context produces moderate degradation, indicating that hypothesis clutter contributes but is not the dominant factor.
- **Panel C:** Varying open goal count (1–8) at a fixed 4096-token context produces mild degradation through reduced goal-focus fidelity.

Within the assumptions of our model, context length is the strongest predictor of failure, but hypothesis clutter and goal count contribute additional degradation. In real agents, these factors are correlated with context length, likely compounding the effect.

## 5 DISCUSSION

*Simulation, not empirical evidence.* We emphasize that our study is a *simulation-based investigation*. The context-degradation model explicitly encodes the hypothesis that longer contexts degrade performance, so finding that longer contexts cause failure in the simulation is expected. The value of the study lies in: (1) quantifying *how much* degradation suffices to explain the observed failure pattern, (2) showing that a simple parametric model produces behavior consistent with Liu et al.’s observations, (3) demonstrating the robustness of the subagent mitigation across model parameters, and (4) generating falsifiable predictions (e.g., critical threshold, calibration gap magnitude) that can be tested with live agents.

*Interaction with lemma complexity.* The key lemma degrades at shorter context lengths compared to generic lemmas, indicating that context length interacts with intrinsic proof difficulty in the model. The alternating-permutation argument requires sustained multi-step reasoning that is especially vulnerable to attention dilution.

*Subagent strategy as mitigation.* The subagent decomposition works by maintaining a hard cap at 2048 tokens throughout proof search, preventing context growth beyond this limit. This keeps the simulated agent in the high-performance regime, validating the approach as a context management strategy.

*Calibration implications.* The simulated calibration gap (confidence remains high as accuracy collapses) suggests that confidence-based self-monitoring would be insufficient for detecting context-induced failure. If this prediction holds in real agents, explicit context-length-aware fallback mechanisms would be needed.

*Paths to validation.* The most impactful next step would be fitting the model parameters ( $L_{crit}$ , decay rates, stall thresholds) to logs from actual Lean-agent runs. Even 20–50 real runs across a few context-length regimes could validate or falsify the predicted threshold and calibration gap. An alternative would be to use the predictions as priors in a Bayesian analysis of sparse real-world data.

*Limitations.* (1) Parameters are chosen to be consistent with published findings but are not fitted to real agent trace data. (2) The model treats context length as a scalar, abstracting over the position and structure of information within the context. (3) The alternative-hypothesis analysis uses simplified models of hypothesis clutter and goal count; real interactions may be more complex. (4) We do not model other sources of difficulty such as library knowledge requirements or type-theoretic complexity.

## 6 CONCLUSION

We have presented a simulation-based investigation of whether long Lean proof contexts cause the observed difficulty of LLM agents on the Putnam 2025 A5 key lemma. Using a calibrated context-degradation model, we find that context length is the strongest

predictor of failure among the factors we modeled, with Spearman  $\rho = -0.85$  and large effect sizes. The subagent decomposition strategy, which hard-caps context at 2048 tokens throughout proof search, robustly mitigates this failure. Sensitivity analysis confirms these conclusions hold across a wide range of model parameters. The calibration gap prediction and the identified critical threshold provide concrete targets for live-agent validation. Our framework demonstrates the utility of simulation-based analysis for understanding failure modes in LLM-based theorem provers while clearly delineating the boundary between model-based predictions and empirical findings.

## REFERENCES

- [1] Anthropic. 2024. The Claude Model Family. *Technical Report* (2024).
- [2] Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q Jiang, Jia Deng, Stella Biderman, and Sean Welleck. 2024. Llemma: An Open Language Model For Mathematics. *International Conference on Learning Representations* (2024).
- [3] Leonardo de Moura and Sebastian Ullrich. 2021. The Lean 4 Theorem Prover and Programming Language. In *International Conference on Automated Deduction*. Springer, 625–635.
- [4] Emily First, Markus N Rabe, Talia Ringer, and Yuriy Brun. 2023. Baldur: Whole-Proof Generation and Repair with Large Language Models. *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering* (2023), 1229–1241.
- [5] Albert Qiaochu Jiang, Sean Welleck, Jin Peng Zhou, Timothée Lacroix, Jiacheng Lutfi, Wenda Matber, Manzil Dwiwedi-Yu, Marie-Anne Lachaux, Yin Li, Julien Sablayrolles, et al. 2023. Draft, Sketch, and Prove: Guiding Formal Theorem Provers with Informal Proofs. *International Conference on Learning Representations* (2023).
- [6] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic Uncertainty: Linguistic Invariances for Uncertainty Estimation in Natural Language Generation. In *International Conference on Learning Representations*.
- [7] Guillaume Lample, Marie-Anne Lachaux, Thibaut Lavril, Xavier Martinet, Amaury Hayat, Gabriel Ebner, Aurélien Rodriguez, and Timothée Lacroix. 2022. HyperTree Proof Search for Neural Theorem Proving. *Advances in Neural Information Processing Systems* 35 (2022).
- [8] Mosh Levy, Alon Jacoby, and Yoav Goldberg. 2024. Same Task, More Tokens: the Impact of Input Length on the Reasoning Performance of Large Language Models. *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics* (2024).
- [9] Tianle Li et al. 2024. Long-context LLMs Struggle with Long In-context Learning. *arXiv preprint arXiv:2404.02060* (2024).
- [10] Jia Liu et al. 2026. Numina-Lean-Agent: An Open and General Agentic Reasoning System for Formal Mathematics. *arXiv preprint arXiv:2601.14027* (2026).
- [11] Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranajape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the Middle: How Language Models Use Long Contexts. *Transactions of the Association for Computational Linguistics* 12 (2024), 157–173.
- [12] Stanislas Polu and Ilya Sutskever. 2020. Generative Language Modeling for Automated Theorem Proving. *arXiv preprint arXiv:2009.03393* (2020).
- [13] Ofir Press, Noah A Smith, and Mike Lewis. 2022. Train Short, Test Long: Attention with Linear Biases Enables Input Length Generalization. *International Conference on Learning Representations* (2022).
- [14] Zijian Wu et al. 2025. InternLM2.5-StepProver: Advancing Automated Theorem Proving via Expert Iteration on Large-Scale LEAN Problems. *arXiv preprint arXiv:2410.15700* (2025).
- [15] Huajian Xin et al. 2024. DeepSeek-Prover: Advancing Theorem Proving in LLMs through Large-Scale Synthetic Data. *arXiv preprint arXiv:2405.14333* (2024).
- [16] Kaiyu Yang, Aidan M Swope, Alex Gu, Rahul Chalapathi, Peiyang Song, Shixing Yu, Maruan Al-Shedivat, Jian Lei, Pengfei Xia, Rui Qin, et al. 2024. LeanDojo: Theorem Proving with Retrieval-Augmented Language Models. *Advances in Neural Information Processing Systems* 36 (2024).