

From Proof to Program: Characterizing Tool-Induced Reasoning Hallucinations in Large Language Models

Farima Fatahi Bayat, Pouya Pezeshkpour, Estevam Hruschka

Megagon Labs

{farima, pouya, estevam}@megagon.ai

Abstract

Tool-augmented Language Models (TaLMs) can invoke external tools to solve problems beyond their parametric capacity. However, it remains unclear whether these tool-enabled gains reflect trustworthy reasoning. Focusing on the Code Interpreter tool, we show that even when tools are selected and executed correctly, TaLMs treat tool outputs as substitutes for reasoning, producing solutions that appear correct but lack coherent justification. We term this failure mode **Tool-Induced Myopia (TIM)**, and study it using PYMATH, a benchmark of 1,679 competition-level mathematical problems for which Python code is *helpful but not sufficient*. We further develop a multi-dimensional evaluation suite to quantify reasoning degradation in TaLMs relative to their non-tool counterparts. Our findings reveal that while TaLMs achieve up to a 19.3 *percentage point* gain in final-answer accuracy, their reasoning behavior consistently deteriorates (e.g., non-tool LLMs win up to 41.5% more often in pairwise comparisons of reasoning process). This degradation intensifies with tool use; the more frequently a model invokes tools, the less coherent its reasoning becomes. Moreover, tool use shifts errors from arithmetic mistakes toward global reasoning failures (logic, assumption, creativity); with TIM present in ~55% of high-risk cases. Finally, we propose a preference-optimization-based framework that realigns TaLMs to use tools as assistive evidence, improving both final-answer accuracy and reasoning depth under tool use. Codes and data are available at: <https://github.com/megagonlabs/TIM>.

1 Introduction

Large Language Models (LLMs) have grown increasingly capable, yet relying solely on their parametric knowledge introduces key limitations, including the inability to access real-time or domain-specific information (Yu and Ji, 2024; Wang et al., 2025), perform precise computations (Lu et al.,

2023b), or fully comprehend user intentions (Qian et al., 2024). To address these shortcomings, Tool-augmented Reasoning (TIR) (Schick et al., 2023; Gou et al., 2024) has emerged as a promising paradigm. It enables LLMs to integrate natural language reasoning with external tools. A tool, in this context (Wang et al., 2024b), is a function interface to an external computer program, where the model generates function calls to interact with it. Frontier LLMs (OpenAI, 2025b; GoogleAI, 2025c; Anthropic, 2025a) now offer native, sandboxed execution for select tools, choosing when to call, executing, and integrating results.

While tool calling significantly extends LLMs’ utility in computation, retrieval, and procedural tasks, it also introduces new sources of failures. The most fundamental errors arise from *tool hallucinations*, where models either select inappropriate tools or misuse them, leading to incorrect or irrelevant outputs (Patil et al., 2025; Xu et al., 2025). Yet even with correct tool selection and successful execution, LLMs can produce non-factual outputs or flawed reasoning. Prior work ties factual errors to conflicts between parametric and retrieved knowledge (Sun et al., 2025) and to the propagation of errors from retrieved content (Magesh et al., 2024). Additional studies show that LLMs’ unrestricted access to external tools can induce tool overuse (Qian et al., 2025), which confuses the model and harms performance. It also encourages cognitive offloading (Wang et al., 2025), which limits the model’s use of its internal reasoning capabilities. Despite these observations, it remains unclear whether tool-augmented reasoning is hallucination-free even when other failure modes are controlled.

In this work, we exclusively focus on the Code Interpreter tool to ensure that tools are both invoked and executed correctly. This setup allows for strictly controlled execution conditions, including correct tool invocation, error-free runs, intended outputs, while avoiding confounds from multiple

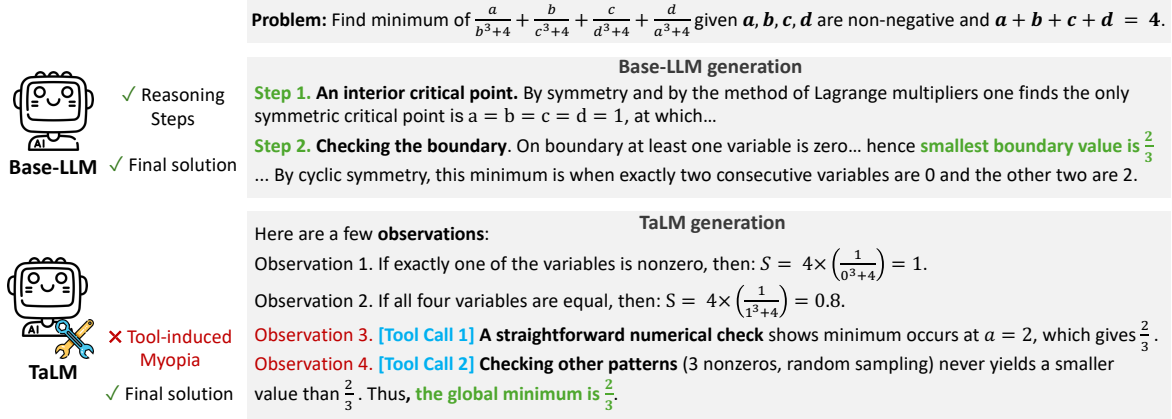


Figure 1: Comparison of Base LLM and Tool-augmented LLM (TaLM) reasoning. The Base LLM (top) derives the solution through step-by-step mathematical reasoning, while the TaLM (bottom) relies on **empirical checks** and multiple tool calls to search for the minimum, a failure mode characteristic of Tool-Induced Myopia (TIM).

tools (e.g., API failures, interface mismatches, retrieval drift, or external data quality issues (Zhong et al., 2025; Faghih et al., 2025; Maekawa et al., 2025; Han et al., 2025)) that would undermine a controlled investigation. We demonstrate that a new family of hallucinations still emerges under these idealized conditions. We term this type of hallucination **Tool-Induced Myopia (TIM)**: a failure mode in which access to an external tool (e.g., a Code Interpreter) causes the model to narrow its reasoning to what the tool can compute, rather than utilizing its full internal reasoning abilities.

Figure 1 illustrates this behavior: with Code Interpreter access (TaLM response), the model repeatedly performs **empirical checks** instead of producing the required reasoning steps generated by the same model without tools (Base-LLM response). The tool is used correctly and returns valid outputs, yet the model’s reasoning depth diminishes. Importantly, existing evaluation approaches cannot capture this failure. Final-answer accuracy cannot detect it, since both responses are correct, and step-level logical-consistency metrics (Lightman et al., 2023; Zheng et al., 2025; Xia et al., 2025; Liu and Fang, 2025) also fail, because the reasoning appears coherent despite skipping essential logical steps. Consequently, exposing TIM requires a richer, multi-dimensional evaluation that examines how the model reasons under tool use.

TIM has significant real-world implications. Models achieving correct answers through opaque reasoning appear reliable but are unsafe in practice; they mislead users about actual reasoning capabilities and erode trust in deployed systems.

To investigate TIM, we introduce PYMATH, a dataset comprising 1,679 competition-level mathe-

matical problems collected from multiple sources and curated to elicit, measure, and mitigate TIM hallucination. We specifically target problems for which code-based computation is helpful but not sufficient for a complete solution (see Section 3.1 for details). We then quantify reasoning degradation in TaLMs compared to their Base (non-tool) counterparts using a comprehensive suite of reference-free and reference-based metrics (Section 3.2). Our results show that, even when TaLMs produce correct final answers (up to a 19.3 *percentage point* in final performance gains), they often overrely on tool outputs and produce shallower reasoning chains than non-tool models. We observed that this degradation intensifies with tool use: the more often a model calls tools, the less coherent its reasoning becomes. Errors also shift from arithmetic mistakes to global reasoning failures (logic, assumption, creativity). Moreover, in a manual audit, we find TIM in ~55% of high-risk model-generated solutions.

To mitigate this issue, we develop two complementary strategies:

To mitigate this issue, we develop two complementary strategies. First, we propose a prompting intervention that encourages models to treat tools as reasoning aids, recovering reasoning behavior at inference time without retraining. Second, we develop a preference-optimization framework based on Direct Preference Optimization (DPO; Rafailov et al., 2023) that aligns TaLMs to integrate the Code Interpreter as an assistant that supports, rather than replaces, mathematical reasoning. On the evaluation split of PYMATH, the fine-tuned TaLM demonstrates improved reasoning behavior and surpasses both the vanilla TaLM (+0.6%) and the Base LLM

(+3.0%) in final-answer accuracy.

In summary, our key contributions are:

- **Tool-induced Myopia:** We uncover a new class of TaLM hallucinations called TIM. We demonstrate that, even under ideal tool-use conditions, LLMs exhibit a systematic failure mode in which tool access suppresses internal reasoning and induces tool-driven shortcuts.
- **Evaluation Benchmark:** To surface TIM, we introduce PYMATH, a dataset of competition-level mathematical problems with step-by-step open-ended solutions, where coding can assist in problem solving. We then provide a multi-dimensional evaluation suite that comprehensively measures the reasoning behavior of TaLM compared to its no-tool counterpart. Our evaluation yields actionable insights into when and why TIM emerges and guides safer tool integration.
- **Mitigation Strategies:** We propose two complementary mitigation approaches: (i) a training-free prompting method and (ii) a preference-optimization-based fine-tuning regime, both reduce TIM and improve reasoning depth in TaLMs.

2 Tool-induced Myopia

We define **Tool-Induced Myopia (TIM)** as a type of hallucination where access to a tool (e.g., Code Interpreter) narrows the model’s reasoning to what the tool can compute, which may result in a shift from step-by-step reasoning to exploratory, tool-driven search. In practice, it substitutes enumeration for proof, skips necessary derivations, mistakes empirical checks for universal guarantees (e.g., brute-forcing search), and may prematurely stop once code returns a plausible output. Note that using code solely for precise computation (e.g., evaluating a determinant or numerically finding a root), while also providing the necessary derivations, is **not** considered TIM. Figure 1 illustrates this hallucination: the Base-LLM solves the problem through step-by-step mathematical reasoning, whereas the TaLM shifts to an exploration-based approach, making empirical observations—some computed internally and others via the Code Interpreter—and arrives at the correct answer through numerical search. TIM degrades the reasoning behavior by overrelying on tools instead of using their outputs as helpful hints. Notably, the TaLM

in Figure 1 still produces the correct final answer via a logically coherent but incomplete sequence of steps. Therefore, evaluations that rely solely on final-answer accuracy or step-level logical consistency measures would overlook the flawed reasoning process and inflate LLM’s performance. To surface this failure mode, we introduce an evaluation benchmark accompanied by a suite of evaluation metrics in what follows.

3 Evaluating TIM in Language Models

3.1 PYMATH

We focus on the domain of mathematical problem solving, where reasoning and computation are tightly coupled, and evaluate LLMs under two settings: with and without access to a Python Code Interpreter. We collect text-only math problems in English from multiple competition-level sources. Table 1 reports the data sources and the number of problems drawn from each dataset. To mitigate potential data contamination, we restrict AIME to 2024–2025 problems. From Omni-Math (Gao et al., 2024), which provides a 1–10 difficulty rating, we retain only problems with difficulty ≥ 5 , since frontier LLMs have already saturated performance on easier problems (OpenAI, 2025a; Anthropic, 2025b; MetaAI, 2025).

To more effectively elicit the TIM hallucination in TaLMs, we target problems for which Python is *helpful but not sufficient*. This setup creates a natural opportunity for tool use. It allows the model to benefit from the Code Interpreter, while still *requiring* reasoning on top of tool outputs to derive a solution. To identify such instances, we adopt an LLM-as-a-judge protocol (Gu et al., 2025) to assess each problem according to two criteria: **(1) Python Usefulness:** whether using Python code helps solve the problem; and **(2) Python Sufficiency:** whether Python code alone (without additional LLM reasoning) is sufficient to fully solve the problem (assessment prompt in Appendix A.1).

Table 1 summarizes the resulting statistics by source. Our final dataset comprises 1,679 competition-level problems with step-by-step reference solutions: 1,000 problems (~60%, distributed across sources and problem difficulties) for evaluation, and the rest for TaLM training (10% development, 90% DPO fine-tuning; detailed in Section 6).

Source	Total	Useful & Not Sufficient	Eval Split (%)
AIME* (AIM, 2024-2025)	60	13	23.1
OlympiadBench (He et al., 2024)	674	171	52.0
OlympicArena (Huang et al., 2024)	169	47	12.8
Omni-Math (Gao et al., 2024)	2569	1448	62.3
PYMATH	-	1679	59.5

Table 1: Sources and statistics of competition-level math problems included in the PYMATH benchmark. *Only AIME 2024–2025 problems included.

3.2 Evaluation Suite

Evaluating LLMs’ mathematical reasoning has traditionally relied on outcome-based evaluation, most predominantly final-answer accuracy (Liu et al., 2025a; Ahn et al., 2024). However, recent studies (Mondorf and Plank, 2024; Yee et al., 2024) have shown that LLMs can reach a correct final answer despite invalid reasoning. This gap has motivated a shift toward process-based evaluation, which assesses LLMs’ *reasoning behavior* rather than their task performance.

Process-based approaches span two major families. *Reference-free* methods do not rely on gold reasoning traces and instead assess reasoning behavior via mechanisms such as self-consistency (Liu and Fang, 2025), or pairwise comparison (e.g., win rate, (Chen et al., 2025; Qin et al., 2024b)). In contrast, *reference-dependent* approaches evaluate intermediate reasoning against a gold step-by-step solution, enabling fine-grained detection of invalid, skipped, or inconsistent steps (Yan et al., 2025; Chernyshev et al., 2025). A recent line of work extends this idea via process reward models (PRMs) to score partial solutions step by step, providing a learned measure of reasoning soundness without requiring gold supervision (Lightman et al., 2023; Zheng et al., 2025; Li et al., 2025).

As illustrated in Figure 1, neither final-answer accuracy nor unidimensional process-based metrics (e.g., PRM scores alone) are sufficient to expose the TIM in TaLMs: the model may output the correct answer and generate seemingly coherent reasoning steps while still exhibiting degraded reasoning under tool use. This motivates the need for a multi-dimensional evaluation suite that jointly captures: (i) task outcome, (ii) counterfactual impact of tool use, (iii) divergence from ground truth reasoning traces, and (iv) step-level logical consistency. Next, we introduce these four evaluation dimensions and clarify why each one is necessary, yet insufficient on its own, for diagnosing TIM.

3.2.1 Final-Answer Accuracy

We first measure final-answer accuracy as the standard measure of task success. Although this metric cannot reveal whether a TaLM solved the problem through mathematical reasoning or tool-driven shortcuts, it confirms that TIM affects reasoning, not task performance. For all subsequent metrics, which target reasoning behavior rather than outcome correctness, we evaluate models only on problems where the final answer is correct. This allows us to isolate the effect of tool use on how the model reasons, rather than whether it succeeds.

3.2.2 Win Rate

To assess whether tool use meaningfully affects reasoning behavior, we compare solutions from a TaLM to those of the corresponding Base (non-tool) LLM. Following prior work demonstrating that LLMs can reliably approximate human preferences (Zheng et al., 2023; Chen et al., 2025), we employ an LLM judge to evaluate reasoning behavior in each response using a rubric adapted from Petrov et al. (2025), covering: (1) *logic errors* (logical fallacies or unjustified leaps), (2) *assumption errors* (unsupported or incorrect assumptions), (3) *creativity errors* (invalid solution strategies), and (4) *algebra/arithmetic errors* (critical symbolic or numeric mistakes). The judge then selects the response with higher reasoning depth and fewer errors as the winner (evaluation prompt in Appendix A.2).

Win Rate measures the fraction of comparisons in which the TaLM solution is preferred (and vice versa). A decline in Win Rate indicates that tool access *harms* the LLM’s reasoning behavior. Notably, this comparison-based setup reduces single-judge bias (Liu et al., 2025b; Fatahi Bayat et al., 2025) and measures the counterfactual effect of tool use.

3.2.3 Miss Rate

Inspired by recall-based measures in long-form factuality evaluation (Wei et al., 2024; Liu et al., 2025b), we define the *Miss Rate* as the proportion of reasoning steps in the ground-truth solution that are absent from the LLM-generated solution, i.e., $\frac{|\text{missing steps}|}{|\text{gold steps}|}$. A high Miss Rate indicates that the model abandoned a valid derivation, skipped necessary steps, or replaced reasoning with trial-and-error code execution. Unlike final-answer accuracy, the Miss Rate detects invalid reasoning paths even when the model reaches the correct solution. However, this metric may over-penalize a solution when

the reference solution(s) do not cover all valid reasoning paths (prompt in Appendix A.3).

3.2.4 PRM Accuracy

Finally, we evaluate step-level reasoning behavior using a Process Reward Model, specifically QWEN2.5-MATH-7B-PRM800K (Zheng et al., 2025), trained on $\sim 800K$ annotated mathematical reasoning steps to identify erroneous ones. An overall correctness score for a solution is obtained by aggregating these step-level scores. While PRMs assess step-level correctness, they may fail to reliably detect reasoning shortcuts where code outputs substitute for derivation, as they do not capture a holistic view of solution completeness.

These four metrics jointly provide minimal yet complete coverage of the TIM construct: Final-answer accuracy captures outcome; Win Rate isolates the tool’s counterfactual impact on reasoning behavior; Miss Rate quantifies divergence from valid reference reasoning; and a PRM assesses step-level soundness without relying on the ground truth. This set is sufficient because TIM exhibits a clear and convergent pattern: final-answer accuracy stays the same or improves in TaLMs, while Win Rate declines, Miss Rate increases, and PRM accuracy drops. This suite both distinguishes TIM from alternative explanations (e.g., general model weakness or genuinely helpful tool use) and is falsifiable: when tools are removed, the pattern recedes, and when tool reliance is increased, it strengthens. Together, these properties establish both the existence and mechanism of Tool-Induced Myopia.

4 Experimental Setup

Large Language Models: We evaluate proprietary frontier language models equipped with in-house Code Interpreters. These models can autonomously decide when to invoke the interpreter, allowing us to study tool-augmented reasoning without modifying or intervening in their internal execution pipeline. We benchmark seven models across three major LLM families: (1) OpenAI models: GPT-4.1-mini and GPT-4.1 (OpenAI, 2025c) (non-thinking LLMs), and o4-mini (OpenAI, 2025d) and GPT-5 (OpenAI, 2025a) (thinking LLMs), (2) Gemini models: Gemini-2.0-Flash (GoogleAI, 2025a) (non-thinking) and Gemini-2.5-Flash (GoogleAI, 2025b) (thinking), (3) Claude-Opus-4 (Anthropic, 2025b) (thinking).

Evaluation Benchmark: We use the evaluation split of PYMATH as our benchmark dataset and

apply our four-dimensional evaluation suite to measure reasoning and tool-use behavior. The suite comprises: (1) Final-answer Accuracy, (2) Win Rate (relative correctness vs. the Base-LLM), (3) Miss Rate (proportion of missing steps from the ground truth), (4) PRM-Evaluated Accuracy (step-level reasoning correctness). For the first three metrics, we adopt an LLM-as-a-judge protocol with GPT-5 serving as the judge. This setup enables a comprehensive comparison between TaLMs and their non-tool (Base) counterparts.

5 Results and Analyses

In this section, we first evaluate Base-LLMs and their tool-augmented counterparts (TaLMs) using our four-dimensional evaluation suite in Section 5.1. Section 5.2 further examines how the severity of TIM in the LLM-generated solution changes as the number of tool calls increases. We also investigate whether the complexity of TaLM-generated code correlates with the degree of reasoning hallucination in their solutions in Section 5.3. Next, we apply the error taxonomy introduced in Section 3.2.2 to analyze the prevalent error types that occur when LLMs have access to a Code Interpreter compared to Base-LLMs (Section 5.4). Since PYMATH encourages but does not guarantee tool use, in Section 5.5 we analyze the frequency of tool invocation for TaLMs. Finally, Section 5.6 presents a qualitative analysis validating that the high-risk solutions flagged by our evaluation suite indeed exhibit TIM behaviors.

5.1 Base-LLMs Show Stronger Reasoning Despite Lower Accuracy

Our benchmarking results on the evaluation set of PYMATH are presented in Table 2. We report Final-answer Accuracy, Miss Rate, Win Rate, and PRM Accuracy across the Base and TaLM variants of seven LLMs. The results reveal a consistent pattern confirming the presence of TIM: although **TaLMs achieve higher Final-answer Accuracy, Base-LLMs exhibit stronger mathematical reasoning**, reflected in lower Miss Rates, higher Win Rates, and higher PRM Accuracy on average. Prior work has shown that LLMs can produce the correct final answer despite flawed or incomplete reasoning (Mondorf and Plank, 2024; Lightman et al., 2023). Our results demonstrate that access to external tools often amplifies this discrepancy.

Interestingly, the highest PRM scores are

Model	Variant	Final Acc. (\uparrow)	Miss Rate (\downarrow)	Win Rate (\uparrow)	PRM Acc. (\uparrow)
GPT-4.1-mini	Base	30.0	45.7	58.6	93.0
	TaLM	28.7	49.9	41.4	88.9
GPT-4.1	Base	24.6	48.1	54.4	88.6
	TaLM	27.0	49.9	45.6	85.9
o4-mini	Base	45.1	45.5	49.0	73.6
	TaLM	64.4	47.6	50.9	67.9
GPT-5-Thinking	Base	67.5	38.8	56.0	57.5
	TaLM	71.9	43.8	44.0	50.2
Gemini-2.0-Flash	Base	24.3	54.2	52.7	65.0
	TaLM	25.1	56.6	47.3	68.5
Gemini-2.5-Flash	Base	45.4	40.2	54.6	81.5
	TaLM	45.7	40.9	45.4	78.8
Claude-Opus-4	Base	28.0	50.9	41.3	77.9
	TaLM	40.5	52.8	58.6	57.8
Average	Base	37.6	45.9	52.4	76.7
	TaLM	43.3	48.8	47.6	71.1

Table 2: Performance of Base-LLMs and TaLMs across four reasoning metrics. The best score in each column is highlighted in green. On average, TaLMs achieve higher Final-answer Accuracy (Final Acc.), but Base-LLMs exhibit greater reasoning depth—as indicated by lower Miss Rate, higher Win Rate, and higher PRM Accuracy—confirming the presence of TIM.

achieved by non-thinking GPT-4.1 models compared to their stronger “thinking” variants. This aligns with recent findings that step-level reward models struggle to reliably assess long and complex reasoning chains, as they often conflate fluency with correctness, become miscalibrated on stronger models, and generalize poorly to the longer, self-correcting traces produced by large reasoning models (Bamba et al., 2025; Lee et al., 2025).

Finally, we posit that the overall gap between Base-LLM and TaLM performance appears moderate in aggregate because TIM is partially masked by three factors: (1) our filtering for TIM-prone problems, while targeted, is not perfect; (2) existing metrics have limited sensitivity to TIM hallucinations; and (3) tool use yields only modest accuracy gains on PYMATH for most LLMs. However, we next show that the reasoning gap widens substantially as tool invocations increase.

5.2 Base-TaLM Reasoning Gap Widens with Tool Call Frequency

We examine how reasoning behavior changes with tool-call frequency, a key driver of TIM: as reliance on tools grows, models increasingly replace reasoning with computation. We test whether higher call counts yield larger Base-TaLM divergence (i.e., higher Miss Rates, Lower Win Rates, and PRM Ac-

curacies) even when final answers are correct. We group problems into bins based on number of tool calls in their LLM-generated solutions: {0–3, 4–7, 8–11, 12+}. As shown in Figure 2, three consistent patterns emerge across most models. First, the *Win Rate* of the Base model against TaLM increases with tool call frequency, indicating that reasoning weakens as solutions rely more heavily on tools. Second, the *Miss Rate* generally increases with more tool usage, reflecting larger deviations from mathematically grounded reasoning. Third, *PRM Accuracy* typically decreases as tool calls grow, suggesting that longer, tool-heavy trajectories accumulate more step-level errors. Occasional reversals in the 12+ bin are due to the limited number of samples in that category.

5.3 Code Complexity Does Not Explain TIM

If TIM were driven by the complexity of generated code, we would expect strong correlations between code complexity and reasoning behavior. We investigate this hypothesis by measuring whether the complexity of TaLM-generated code correlates with TIM severity. We measure code complexity using two standard metrics (Dou et al., 2024; Chen et al., 2024): (1) *Line of Code* (Boehm, 1981) (average number of lines per code block), and (2) *Cyclomatic Complexity* (McCabe, 1976),

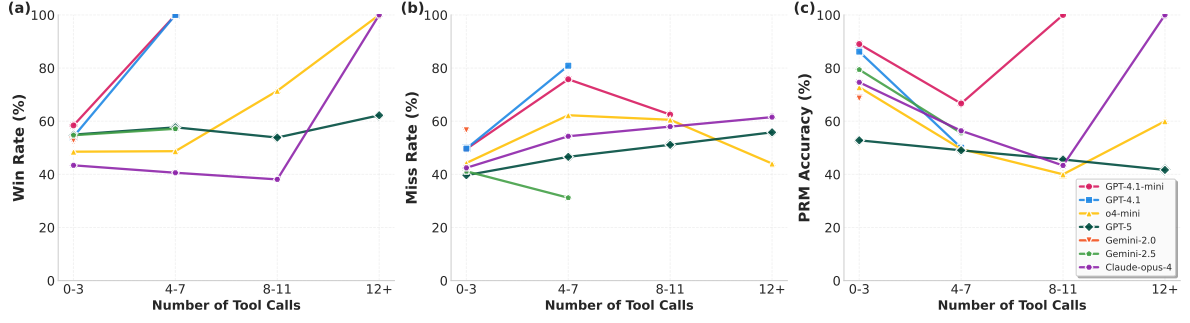


Figure 2: **Reasoning behavior vs. tool usage.** Metrics are computed over bins defined by the number of tool calls in model solution: {0–3, 4–7, 8–11, 12+}. **(a) Win Rate** (higher is better) shows that the Base model more frequently outperforms TaLM as tool usage increases, indicating a widening gap at higher call counts. **(b) Miss Rate** (higher is worse) generally rises with additional tool calls, indicating more missing steps in tool-dependent trajectories. **(c) PRM Accuracy** (higher is better) typically declines as the number of calls grows.

which measures the number of linearly independent paths through a program’s source code. We compute Pearson correlations between these complexity metrics and Miss Rate, which provides continuous values suitable for correlation analysis. Figure 3 shows correlation coefficients (left) and their statistical significance (right) across all TaLMs. **Overall, we find no statistically significant correlation between Miss Rate and either complexity metric:** Line of Code correlations range from -0.09 to 0.22, while Cyclomatic Complexity correlations range from -0.03 to 0.26. Although a few models show marginal correlations at the $p < 0.10$ level, these effects do not reach statistical significance and are inconsistent across models. Together, these results indicate that the syntactic or structural complexity of generated code does not drive TIM.

5.4 Tool Use Shifts Errors from Arithmetic to Global Reasoning Failures

In this section, we analyze shifts in reasoning error types by comparing TaLM solutions with their Base-LLM counterparts. Our goal is to understand TIM’s core mechanism: how reasoning changes when a model shifts from mathematical derivation to computation-heavy solutions. We use the error taxonomy from Section 3.2.2 (inspired by Petrov et al. (2025)) and prompt our judge (GPT-5) to rate the presence of each error type on a 1–5 scale, or output None if no error is detected (prompt in Appendix A.4). To isolate changes in reasoning behavior, we only analyze cases where **both** Base-LLM and TaLM produce the correct final answer.

Across models, we observe a consistent shift in error patterns with tool use (Figure 4): logic, assumption, and creativity errors in almost all cases increase, indicating that TaLMs make more unjusti-

fied inferences and reasoning leaps when relying on tool outputs. In contrast, algebraic and arithmetic errors decrease, since computation is delegated to Code Interpreter, while the proportion of “no error” cases drops, showing that tool access introduces new reasoning flaws even when final answers remain correct. Therefore, TIM does not simply result in missing steps; rather, it fundamentally reshapes the model’s *reasoning behavior*.

5.5 Thinking Models Use Code More Frequently

Having shown that more tool calls intensify TIM (which in turn widens Base–TaLM divergence and degrades reasoning metrics), we now quantify how often each model invokes the Code Interpreter tool on math problems, providing a “dose” measure that clarifies TIM risk across models. We investigate tool invocation frequencies across TaLMs on PY-MATH, which was curated to encourage code use by selecting problems where computation is helpful. Figure 5 shows the percentage of problems on which each TaLM invoked the Code Interpreter tool. Thinking models exhibit substantially higher tool use rates compared to their non-thinking counterparts: Claude-opus-4 achieves the highest rate at 99.8%, followed by GPT-5 at 73.7%, while non-thinking models show more modest usage—GPT-4.1-mini (16.4%), GPT-4.1 (25.9%), and Gemini-2.0-flash (14.3%). **On average, thinking models invoke tool on 49.7% more problems than non-thinking models.** Given the established relationship between tool call frequency and TIM severity (Section 5.2), this heavy reliance on computation suggests thinking models face higher TIM risk.

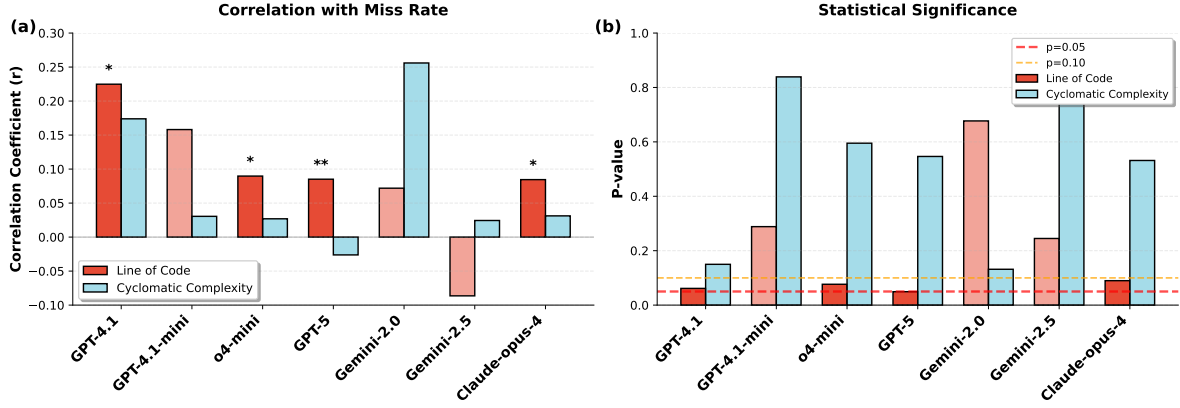


Figure 3: Correlation between code complexity metrics and Miss Rate across TaLMs. (a) Pearson correlation coefficients for Line of Code and Cyclomatic Complexity with Miss Rate. (b) Corresponding p-values with significance thresholds at $p=0.05$ (red) and $p=0.10$ (orange). Asterisks denote marginal significance (* $p<0.10$, ** $p<0.05$). No statistically significant correlations are found, suggesting that TIM is not driven by code complexity.



Figure 4: Change in reasoning error rates after tool use (Δ = TaLM – Base). Positive values indicate that an error type becomes more frequent when the model has access to Code Interpreter tool.

5.6 TIM in Over Half of High-Risk Solutions

We conduct a targeted manual evaluation to validate TIM qualitatively and to surface recurrent linguistic cues that precede it. We focus on high-risk LLM-generated solutions where: (1) the final answer is correct, (2) the Base-LLM outperforms TaLM in Win Rate judgment, and (3) the PRM flags errors in TaLM’s reasoning traces. From this filtered set, we select the top-10 samples per model ranked by Miss Rate and examine them to determine whether TIM is present and identify characteristic precursor phrases.

Our inspection reveals that **TIM appears in 54.3% of high-risk cases across models**, demonstrating the correlation between TIM incidence and our automated metrics. TIM prevalence varies across models: those where tool use has a greater impact on reasoning quality (Table 2), such as Claude-Opus-4 and o4-mini, exhibit more frequent TIM, while Gemini models, for which tool use has

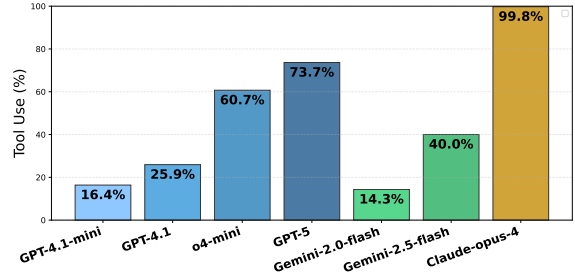


Figure 5: Code Interpreter invocation rates across TaLMs. Thinking models (GPT-5, o4-mini, Gemini-2.5-Flash) use the tool on an average of ~50% more problems than non-thinking models.

minimal impact, show lower rates. Notably, less capable models often display explicit linguistic cues immediately before TIM instances, such as “one numerically finds,” “systematic checks show,” or “let’s verify [an assumption] programmatically.” In contrast, more capable models like GPT-5 manifest TIM more subtly; they silently substitute code outputs for missing derivations. We provide a comprehensive list of high-frequency precursor phrases in Appendix A.5. These phrases serve as indicators rather than definitive proof of TIM and should be interpreted alongside our quantitative metrics.

6 Mitigating TIM Hallucination

Our evaluation results reveal widespread TIM in TaLM-generated solutions, where models rely excessively on tool outputs rather than mathematical reasoning. In this section, we propose two complementary mitigation strategies that encourage models to use the Code Interpreter as an *assistant to reasoning* rather than a substitute for it.

6.1 Prompting-Based Mitigation

We design a lightweight prompting strategy that self-instructs the model to use tool outputs as reasoning aids. Specifically, after each problem statement, we inject the following instruction as a model-generated message:

“We should treat code snippets and their execution results only as helpful hints, and derive the solution through mathematical reasoning.”

This single-sentence intervention encourages the TaLM to treat tool outputs as *verification aids*, thereby promoting mathematical thinking over computational shortcuts.

6.2 Alignment-Based Mitigation

We develop a training-based framework to mitigate TIM through direct preference optimization (DPO; (Rafailov et al., 2023)). Due to resource constraints, we focus on fine-tuning a single model (GPT-4.1) to demonstrate the effectiveness of this approach.

6.2.1 Preference Data Creation

We construct a preference dataset for DPO training using the training split of PYMATH. For each problem, we generate a *chosen-rejected* pair, where the chosen solution demonstrates high-quality reasoning and the rejected solution exhibits TIM characteristics.

Chosen samples: We generate chosen solutions by applying the prompting strategy from Section 6.1 to GPT-4.1, which encourages balanced reasoning and tool use.

Rejected samples: We generate rejected examples through controlled degradation: given a problem and its corresponding chosen solution (containing interleaved reasoning, code tool calls, and execution outputs), we prompt the same model to rewrite a text span of the solution with explicit excessive reliance on tool outputs. Specifically, the model is instructed to produce a coherent but tool-dependent version that omits or abbreviates intermediate mathematical reasoning steps. This process creates naturalistic rejected samples that reflect the TIM hallucination. More details on dataset creation and experimental setup are provided in Appendix A.6.

6.3 Results

Table 3 presents mitigation results over our evaluation benchmark. Both mitigation strategies reduce

Variant	Final Acc. (↑)	Miss Rate (↓)	Win Rate (↑)	PRM Acc. (↑)
Base	24.6	48.1	54.4	88.6
TaLM	27.0	49.9	45.6	85.9
TaLM + Prompting	25.1	49.4	52.7	82.9
TaLM + DPO	27.6	46.6	58.2	83.3

Table 3: Impact of mitigation strategies on tool-augmented GPT-4.1. Prompting reduces TIM without retraining but at the cost of accuracy, while DPO alignment improves both final answer accuracy and reasoning quality (as measured by Win Rate and Miss Rate).

TIM hallucination compared to the vanilla TaLM, but with distinct trade-offs.

Prompting intervention: The zero-shot prompting strategy substantially recovers reasoning quality without model retraining, improving Win Rate from 45.6% to 52.7%, while Miss Rate decreases slightly from 49.9% to 49.4%. However, this comes at the cost of final-answer accuracy, which drops from 27.0% to 25.1%, closer to the base model’s 24.6%.

DPO alignment: The fine-tuned model achieves strong performance across multiple dimensions: highest final-answer accuracy (27.6%, +0.6% over vanilla TaLM), highest Win Rate (58.2% vs. 54.4% Base), and lowest Miss Rate (46.6%). However, PRM accuracy (83.3%) remains below the base model (88.6%), suggesting that while DPO successfully reduces global reasoning errors and tool over-reliance, it does not fully recover step-level correctness as measured by process reward models. Nonetheless, these results demonstrate that TaLMs can be aligned to leverage tools as *reasoning aids* rather than *reasoning shortcuts*, achieving simultaneous improvements in both correctness and reasoning behavior.

7 Related Work

Tool-augmented reasoning extends LLMs through external programmatic interfaces, enabling them to solve tasks beyond their parametric knowledge. Early work such as Toolformer (Schick et al., 2023) demonstrated that LMs can autonomously learn when and how to call external functions, while subsequent surveys (Wang et al., 2024b; Gou et al., 2024) define LLM-used tools and provide broad overviews of tool-augmented LM capabilities. With the advent of tool-augmented reasoning, numerous tools, including web search (Asai et al., 2024; Lu et al., 2023a), document retrieval systems (Fan et al., 2024), code interpreters (Chen et al., 2023; Gao et al., 2023a), and domain-specific

APIs (Qin et al., 2023), have enhanced the reasoning and problem-solving capacity of modern LMs.

Hallucinations in Tool-Augmented LMs Despite these advances, several failure modes persist in tool-augmented LMs (TaLMs). Prior work has identified multiple sources of hallucinations: *tool hallucinations*, where models either select incorrect tools or misuse them (Xu et al., 2025); inconsistencies caused by conflicts between parametric and non-parametric knowledge (Sun et al., 2025); and factual errors propagated from retrieved content (Magesh et al., 2024). Prior work has also examined attribution failures in retrieval-augmented generation (Gao et al., 2023b; Liu et al., 2023), where models fail to ground their outputs in retrieved evidence. Moreover, unrestricted access to tools can promote *cognitive offloading*, hindering models from developing or utilizing their internal reasoning capabilities (Qian et al., 2025; Wang et al., 2025). However, prior studies have overlooked reasoning degradation that occurs *even when tools are used correctly*. We define this phenomenon as Tool-Induced Myopia (TIM): a failure mode where tool access biases reasoning toward computational shortcuts.

Evaluation Benchmarks in Mathematical Reasoning Evaluating mathematical reasoning in LLMs has traditionally focused on outcome correctness using benchmarks such as GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021). However, these datasets contain problems that frontier LMs have largely mastered (with performance approaching saturation), and outcome-based metrics fail to capture the quality of LLMs’ reasoning processes (Mondorf and Plank, 2024). More recent work has introduced rationale-based evaluation methods, including reward models for step-level validation (Lightman et al., 2023; Wang et al., 2024a; Zheng et al., 2025), intermediate chain-of-thought scoring (Xia et al., 2025; Li et al., 2025), and model-based judges for reasoning quality (Zhou et al., 2025; Xu et al., 2025). Moreover, competition-level benchmarks (AIM, 2024-2025; He et al., 2024; Huang et al., 2024; Gao et al., 2024; Phan et al., 2025; Fang et al., 2024) have raised the difficulty bar. However, these one-dimensional metrics assess only final correctness or surface-level consistency, failing to distinguish robust mathematical reasoning from tool-driven shortcuts. We fill this gap with a benchmark of competition-level math problems designed to elicit TIM, alongside

a multi-dimensional evaluation that measures both the presence and the severity of tool-induced reasoning degradation.

Mitigation Techniques for TaLM Failures Various strategies have been proposed to address failures in tool-augmented systems. Prompting-based approaches include self-reflection (Shinn et al., 2023), where models critique and revise their reasoning. Alignment-based methods have addressed tool selection and tool-calling hallucinations (Qin et al., 2024a; Patil et al., 2024), tool-calling decision-making through preference optimization (Ross et al., 2025), reliability alignment via expanded action spaces (Xu et al., 2025), and hybrid approaches that combine parametric and non-parametric reasoning (Asai et al., 2024). However, these strategies primarily target explicit errors in tool selection and usage rather than subtle reasoning degradation that occurs even when tools are correctly applied. Our work introduces interventions that explicitly reward comprehensive analytical reasoning over premature tool reliance, providing the first systematic approach to mitigating TIM.

8 Conclusion

Investigating tool-augmented reasoning across competition-level math, we showed that even when tools are used correctly, tool-augmented LLMs often exhibit **Tool-Induced Myopia (TIM)**—a failure mode where models substitute external computation for mathematical reasoning. To study this, we introduced PYMATH and a four-dimensional evaluation suite revealing a consistent pattern across seven frontier models: TaLMs improve final-answer accuracy but produce less complete and reliable reasoning traces than their non-tool counterparts. Error analysis confirms that tool use reduces arithmetic mistakes but increases logical and assumption errors, especially as tool calls grow. Finally, we introduced two mitigation strategies: prompting and DPO-based preference optimization, that reduce TIM hallucination and restore trustworthy reasoning.

9 Limitations

While our study provides the first characterization of Tool-Induced Myopia (TIM), it also has several limitations: First, to ensure precise control over tool invocation and execution, we restrict our analysis to a single tool, the Code Interpreter. This choice eliminates confounds such as API failures

or retrieval noise, but it limits the generality of our findings to other tool types (e.g., search, retrieval, or various APIs). Future work should extend this framework to a broader range of tools and interaction settings. Second, our manual investigation of TIM occurrence and linguistic precursors is limited in scope. While qualitative inspection confirms strong alignment with our automated metrics, a larger-scale human evaluation, covering more models, domains, and annotators, would strengthen the robustness and generalizability of these observations. Finally, due to computational and resource constraints, our mitigation experiments focus exclusively on the GPT-4.1. Although this choice provides a representative case for high-end tool-augmented reasoning, reproducing the mitigation analyses on additional LLMs remains an important direction for validating consistency and scalability.

References

- 2024-2025. Aime i and aime ii problems and solutions. https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions.
- Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. 2024. [Large language models for mathematical reasoning: Progresses and challenges](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 225–237, St. Julian’s, Malta. Association for Computational Linguistics.
- Anthropic. 2025a. Claude models. <https://docs.claude.com/en/docs/about-claude/models/overview>. Version: 2025-11-10.
- Anthropic. 2025b. System card: Claude opus 4 and claude sonnet 4. Version: 2025-05.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. [Self-RAG: Learning to retrieve, generate, and critique through self-reflection](#). *Proceedings of the International Conference on Learning Representations (ICLR) 2024*.
- Udbhav Bamba, Heng Yang, Rishabh Tiwari, Michael W. Mahoney, Kurt Keutzer, and Amir Gholami. 2025. [Reward under attack: Evaluating the sensitivity of process reward models](#).
- Barry W. Boehm. 1981. *Software Engineering Economics*. Prentice-Hall, Englewood Cliffs, NJ.
- Liguo Chen, Qi Guo, Hongrui Jia, Zhengran Zeng, Xin Wang, Yijiang Xu, Jian Wu, Yidong Wang, Qing Gao, Jindong Wang, and 1 others. 2024. A survey on evaluating large language models in code generation tasks. *arXiv preprint arXiv:2408.16498*.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.
- Xilun Chen, Ilia Kulikov, Vincent-Pierre Berges, Barlas Oğuz, Rulin Shao, Gargi Ghosh, Jason Weston, and Wen tau Yih. 2025. [Learning to reason for factuality](#). *Preprint*, arXiv:2508.05618.
- Konstantin Chernyshev, Vitaliy Polshkov, Vlad Stepanov, Alex Myasnikov, Ekaterina Artemova, Alexei Miasnikov, and Sergei Tilga. 2025. U-math: A university-level benchmark for evaluating mathematical skills in large language models. In *Proceedings of the Fourth Workshop on Generation, Evaluation and Metrics (GEM²)*, pages 974–1001.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *Preprint*, arXiv:2110.14168.
- Shihan Dou, Haoxiang Jia, Shenxi Wu, Huiyuan Zheng, Weikang Zhou, Muling Wu, Mingxu Chai, Jessica Fan, Caishuang Huang, Yunbo Tao, and 1 others. 2024. What’s wrong with your code generated by large language models? an extensive study. *arXiv preprint arXiv:2407.06153*.
- Kazem Faghieh, Wenxiao Wang, Yize Cheng, Siddhant Bharti, Gaurang Sriramanan, Sriram Balasubramanian, Parsa Hosseini, and Soheil Feizi. 2025. Tool preferences in agentic llms are unreliable. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 20965–20980.
- Wenqi Fan, Yajuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. [A survey on rag meeting llms: Towards retrieval-augmented large language models](#). In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD ’24*, page 6491–6501, New York, NY, USA. Association for Computing Machinery.
- Meng Fang, Xiangpeng Wan, Fei Lu, Fei Xing, and Kai Zou. 2024. [Mathodyssey: Benchmarking mathematical problem-solving skills in large language models using odyssey math data](#). *Preprint*, arXiv:2406.18321.
- Farima Fatahi Bayat, Lechen Zhang, Sheza Munir, and Lu Wang. 2025. [FactBench: A dynamic benchmark for in-the-wild language model factuality evaluation](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 33090–33110, Vienna, Austria. Association for Computational Linguistics.
- Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang

- Chen, Runxin Xu, Zhengyang Tang, Benyou Wang, Daoguang Zan, Shanghaoran Quan, Ge Zhang, Lei Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu, and Baobao Chang. 2024. [Omni-math: A universal olympiad level mathematic benchmark for large language models](#). *Preprint*, arXiv:2410.07985.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023a. Pal: Program-aided language models. *Proceedings of the 40th International Conference on Machine Learning*.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023b. [Enabling large language models to generate text with citations](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6465–6488, Singapore. Association for Computational Linguistics.
- GoogleAI. 2025a. Gemini 2.0: Flash, flash-lite and pro. <https://developers.googleblog.com/en/gemini-2-family-expands/>.
- GoogleAI. 2025b. Gemini 2.5 flash best for fast performance on everyday tasks. <https://deepmind.google/models/gemini/flash/>.
- GoogleAI. 2025c. Gemini models. <https://ai.google.dev/gemini-api/docs/models>. Version: 2025-11-10.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, yelong shen, Yujia Yang, Minlie Huang, Nan Duan, and Weizhu Chen. 2024. [ToRA: A tool-integrated reasoning agent for mathematical problem solving](#). In *The Twelfth International Conference on Learning Representations*.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni, and Jian Guo. 2025. [A survey on llm-as-a-judge](#). *Preprint*, arXiv:2411.15594.
- Ziwen Han, Meher Mankikar, Julian Michael, and Zifan Wang. 2025. Search-time data contamination. *arXiv preprint arXiv:2508.13180*.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. 2024. [OlympiadBench: A challenging benchmark for promoting AGI with olympiad-level bilingual multimodal scientific problems](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3828–3850, Bangkok, Thailand. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. [Measuring mathematical problem solving with the math dataset](#). In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1.
- Zhen Huang, Zengzhi Wang, Shijie Xia, Xuefeng Li, Haoyang Zou, Ruijie Xu, Run-Ze Fan, Lyumanshan Ye, Ethan Chern, Yixin Ye, Yikai Zhang, Yuqing Yang, Ting Wu, Binjie Wang, Shichao Sun, Yang Xiao, Yiyuan Li, Fan Zhou, Steffi Chern, and 9 others. 2024. [Olympicarena: Benchmarking multi-discipline cognitive reasoning for superintelligent ai](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 19209–19253. Curran Associates, Inc.
- Dong Bok Lee, Seanie Lee, Sangwoo Park, Minki Kang, Jinheon Baek, Dongki Kim, Dominik Wagner, Jiongdao Jin, Heejun Lee, Tobias Bocklet, Jinyu Wang, Jingjing Fu, Sung Ju Hwang, Jiang Bian, and Lei Song. 2025. [Rethinking reward models for multi-domain test-time scaling](#). *Preprint*, arXiv:2510.00492.
- Vladimir I. Levenshtein. 1966. [Binary Codes Capable of Correcting Deletions, Insertions and Reversals](#). *Soviet Physics Doklady*, 10(8).
- Ruosun Li, Ziming Luo, and Xinya Du. 2025. [FG-PRM: Fine-grained hallucination detection and mitigation in language model mathematical reasoning](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 4247–4278, Suzhou, China. Association for Computational Linguistics.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. [Let’s verify step by step](#). *Preprint*, arXiv:2305.20050.
- MingShan Liu and Jialing Fang. 2025. [Enhancing mathematical reasoning in large language models with self-consistency-based hallucination detection](#). *Preprint*, arXiv:2504.09440.
- Nelson Liu, Tianyi Zhang, and Percy Liang. 2023. [Evaluating verifiability in generative search engines](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7001–7025, Singapore. Association for Computational Linguistics.
- WenTao Liu, Hanglei Hu, Jie Zhou, Yuyang Ding, Junsong Li, Jiayi Zeng, MengLiang He, Qin Chen, Bo Jiang, Aimin Zhou, and Liang He. 2025a. [Mathematical language models: A survey](#). *ACM Comput. Surv.* Just Accepted.
- Xin Liu, Lechen Zhang, Sheza Munir, Yiyang Gu, and Lu Wang. 2025b. [VeriFact: Enhancing long-form factuality evaluation with refined fact extraction and reference facts](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 17919–17936, Suzhou, China. Association for Computational Linguistics.
- Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2023a. Chameleon: plug-and-play compositional reasoning with large language models. In *Proceedings of the 37th International Conference*

- on *Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA. Curran Associates Inc.
- Pan Lu, Liang Qiu, Wenhao Yu, Sean Welleck, and Kai-Wei Chang. 2023b. [A survey of deep learning for mathematical reasoning](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14605–14631, Toronto, Canada. Association for Computational Linguistics.
- Seiji Maekawa, Jackson Hassell, Pouya Pezeshkpour, Tom Mitchell, and Estevam Hruschka. 2025. Towards reliable benchmarking: A contamination free, controllable evaluation framework for multi-step llm function calling. *arXiv preprint arXiv:2509.26553*.
- Varun Magesh, Faiz Surani, Matthew Dahl, Mirac Suzgun, Christopher D. Manning, and Daniel E. Ho. 2024. [Hallucination-free? assessing the reliability of leading ai legal research tools](#). *Preprint*, arXiv:2405.20362.
- T.J. McCabe. 1976. [A complexity measure](#). *IEEE Transactions on Software Engineering*, SE-2(4):308–320.
- MetaAI. 2025. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation. <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>. Version: 2025-04-05.
- Philipp Mondorf and Barbara Plank. 2024. [Beyond accuracy: Evaluating the reasoning behavior of large language models – a survey](#). In *First Conference on Language Modeling*.
- OpenAI. 2025a. [Gpt-5 system card](#). <https://cdn.openai.com/gpt-5-system-card.pdf>. Version: 2025-08-13.
- OpenAI. 2025b. [Gpt models](#). <https://platform.openai.com/docs/models>. Version: 2025-11-10.
- OpenAI. 2025c. [Introducing gpt-4.1 in the api](#). <https://openai.com/index/gpt-4-1/>. Version: 2025-04-14.
- OpenAI. 2025d. [o4-mini](#). <https://platform.openai.com/docs/models/o4-mini>.
- Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. 2024. [Gorilla: Large language model connected with massive apis](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 126544–126565. Curran Associates, Inc.
- Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. 2025. Gorilla: large language model connected with massive apis. In *Proceedings of the 38th International Conference on Neural Information Processing Systems*, NIPS '24, Red Hook, NY, USA. Curran Associates Inc.
- Ivo Petrov, Jasper Dekoninck, Lyuben Baltadzhiev, Maria Drencheva, Kristian Minchev, Mislav Balunović, Nikola Jovanović, and Martin Vechev. 2025. [Proof or bluff? evaluating llms on 2025 usa math olympiad](#). *Preprint*, arXiv:2503.21934.
- Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, and Mohamed Shaaban et al. 2025. [Humanity's last exam](#). *Preprint*, arXiv:2501.14249.
- Cheng Qian, Emre Can Acikgoz, Hongru Wang, Xiusi Chen, Avirup Sil, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. 2025. [SMART: Self-aware agent for tool overuse mitigation](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 4604–4621, Vienna, Austria. Association for Computational Linguistics.
- Cheng Qian, Bingxiang He, Zhong Zhuang, Jia Deng, Yujia Qin, Xin Cong, Zhong Zhang, Jie Zhou, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2024. [Tell me more! towards implicit user intention understanding of language model driven agents](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1088–1113, Bangkok, Thailand. Association for Computational Linguistics.
- Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Xuanhe Zhou, Yufei Huang, Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su, Huadong Wang, Cheng Qian, Runchu Tian, Kunlun Zhu, Shihao Liang, Xingyu Shen, and 24 others. 2024a. [Tool learning with foundation models](#). *ACM Comput. Surv.*, 57(4).
- Yujia Qin, Shi Liang, Yining Ye, Kunlun Zhu, Lan Yan, Ya-Ting Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Marc H. Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. [Toolllm: Facilitating large language models to master 16000+ real-world apis](#). *ArXiv*, abs/2307.16789.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, dahai li, Zhiyuan Liu, and Maosong Sun. 2024b. [Toolllm: Facilitating large language models to master 16000+ real-world apis](#). In *International Conference on Representation Learning*, volume 2024, pages 9695–9717.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. Direct preference optimization: your language model is secretly a reward model. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA. Curran Associates Inc.
- Hayley Ross, Ameya Sunil Mahabaleshwarkar, and Yoshi Suhara. 2025. [When2Call: When \(not\) to call](#)

- tools. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3391–3409, Albuquerque, New Mexico. Association for Computational Linguistics.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessí, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: language models can teach themselves to use tools. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS ’23*, Red Hook, NY, USA. Curran Associates Inc.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS ’23*, Red Hook, NY, USA. Curran Associates Inc.
- Zhongxiang Sun, Xiaoxue Zang, Kai Zheng, Yang Song, Jun Xu, Xiao Zhang, Weijie Yu, Yang Song, and Han Li. 2025. Redeeep: Detecting hallucination in retrieval-augmented generation via mechanistic interpretability. *Preprint*, arXiv:2410.11414.
- Hongru Wang, Cheng Qian, Wanjuan Zhong, Xiusi Chen, Jiahao Qiu, Shijue Huang, Bowen Jin, Mengdi Wang, Kam-Fai Wong, and Heng Ji. 2025. Acting less is reasoning more! teaching model to act efficiently. *Preprint*, arXiv:2504.14870.
- Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. 2024a. Math-shepherd: Verify and reinforce LLMs step-by-step without human annotations. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9426–9439, Bangkok, Thailand. Association for Computational Linguistics.
- Zhiruo Wang, Zhoujun Cheng, Hao Zhu, Daniel Fried, and Graham Neubig. 2024b. What are tools anyway? a survey from the language model perspective. *Preprint*, arXiv:2403.15452.
- Jerry Wei, Chengrun Yang, Xinying Song, Yifeng Lu, Nathan Hu, Jie Huang, Dustin Tran, Daiyi Peng, Ruibo Liu, Da Huang, Cosmo Du, and Quoc V. Le. 2024. Long-form factuality in large language models. In *Proceedings of the 38th International Conference on Neural Information Processing Systems, NeurIPS ’24*, Red Hook, NY, USA. Curran Associates Inc.
- Shijie Xia, Xuefeng Li, Yixin Liu, Tongshuang Wu, and Pengfei Liu. 2025. Evaluating mathematical reasoning beyond accuracy. In *Proceedings of the Thirty-Ninth AAAI Conference on Artificial Intelligence and Thirty-Seventh Conference on Innovative Applications of Artificial Intelligence and Fifteenth Symposium on Educational Advances in Artificial Intelligence, AAAI’25/IAAI’25/EAAI’25*. AAAI Press.
- Hongshen Xu, Zichen Zhu, Lei Pan, Zihan Wang, Su Zhu, Da Ma, Ruisheng Cao, Lu Chen, and Kai Yu. 2025. Reducing tool hallucination via reliability alignment.
- Yuchen Yan, Jin Jiang, Zhenbang Ren, Yijun Li, Xudong Cai, Yang Liu, Xin Xu, Mengdi Zhang, Jian Shao, Yongliang Shen, and 1 others. 2025. Verify-bench: Benchmarking reference-based reward systems for large language models. *arXiv preprint arXiv:2505.15801*.
- Evelyn Yee, Alice Li, Chenyu Tang, Yeon Ho Jung, Ramamohan Paturi, and Leon Bergen. 2024. Faithful and unfaithful error recovery in chain of thought. In *First Conference on Language Modeling*.
- Pengfei Yu and Heng Ji. 2024. Information association for language model updating by mitigating LM-logical discrepancy. In *Proceedings of the 28th Conference on Computational Natural Language Learning*, pages 117–129, Miami, FL, USA. Association for Computational Linguistics.
- Chujie Zheng, Zhenru Zhang, Beichen Zhang, Runji Lin, Keming Lu, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. 2025. ProcessBench: Identifying process errors in mathematical reasoning. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1009–1024, Vienna, Austria. Association for Computational Linguistics.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging Llm-as-a-judge with mt-bench and chatbot arena. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS ’23*, Red Hook, NY, USA. Curran Associates Inc.
- Lucen Zhong, Zhengxiao Du, Xiaohan Zhang, Haiyi Hu, and Jie Tang. 2025. Complexfunctbench: exploring multi-step and constrained function calling under long-context scenario. *arXiv preprint arXiv:2501.10132*.
- Zihao Zhou, Shudong Liu, Maizhen Ning, Wei Liu, Jindong Wang, Derek Wong, Xiaowei Huang, Qiufeng Wang, and Kaizhu Huang. 2025. Is your model really a good math reasoner? evaluating mathematical reasoning with checklist. In *International Conference on Representation Learning*, volume 2025, pages 34238–34281.

A Appendix

A.1 PYMATH Data Curation

The prompt used for filtering mathematical problems is shown in A.1. We use GPT-5 as the LLM-as-a-judge to make a binary classification decision per problem, a setting where LLMs have been shown to perform reliably (Gu et al., 2025).

Problem Annotation Prompt

You are a technical reasoning assistant for mathematical problem solving. Your task is to evaluate mathematical problems used for benchmarking LLMs in terms of:

1. **Python Usefulness:** Whether using Python code is helpful for solving this problem.
2. **Python Sufficiency:** Whether Python code alone (without extra reasoning steps from the target LLM) is sufficient to fully solve this problem.

Use the following evaluation criteria:

- **Mathematical Domain:** What area(s) of mathematics does this problem involve, and how computational versus theoretical is this domain typically?
- **Solution Type:** What kind of answer or result is the problem asking for?
- **Computational Approach:** What computational strategies, if any, could be applied?
- **Problem Scale:** How do size and complexity affect computational feasibility?
- **Verification Needs:** Would solving the problem benefit from computational verification?
- **Techniques Required:** What mathematical insights or methods are necessary, and how much can be automated?

Based on your evaluation, provide:

- **Python Usefulness:** true/false
- **Python Sufficiency:** true/false
- **Recommendation:** One of "Pure Python", "Python + LLM Insight/Reasoning", "Python for Verification", "Python for Exploration", or "Minimal Python Role".

Respond strictly in the following JSON format:

```
{
  "problem": "repeat problem description here",
  "reasoning": {
    "mathematical_domain": "",
    "solution_type": "",
    "computational_approach": "",
    "problem_scale": "",
    "verification_needs": "",
    "techniques_required": ""
  },
  "python_usefulness": true or false,
  "python_sufficiency": true or false,
  "recommendation": "[one of the options above]"
}
```

A.2 Win Rate Evaluation Prompt

The prompt used by GPT-5 for Win Rate evaluation is shown below.

Win Rate Evaluation Prompt

You are an expert mathematician tasked with grading two solutions, "A" and "B", to the same competition-style problem. Evaluate which solution is better by assigning error severity scores (0 = low, 5 = high) for the following categories:

1. **Logic (0-5):** Errors due to logical fallacies or unjustified leaps in reasoning.
2. **Assumption (0-5):** Errors from unproven or incorrect assumptions that undermine subsequent steps.
3. **Creativity (0-5):** Errors from fundamentally incorrect or unoriginal strategies indicating failure to identify the right approach.
4. **Algebra/Arithmetic (0-5):** Errors arising from critical algebraic or arithmetic miscalculations.

Evaluation:

- Provide a brief justification for each score, referencing relevant mathematical concepts or reasoning techniques.
- Compute a final score as the average of the four error categories.
- Select the solution with the lower final score as the better one. If tied, prefer the solution with clearer reasoning.

Respond strictly in the following JSON format:

```
{
  "A_grades": {
    "logic": {"score": 0-5, "explanation": ""},
    "assumption": {"score": 0-5, "explanation": ""},
    "creativity": {"score": 0-5, "explanation": ""},
    "algebra_arithmetic": {"score": 0-5, "explanation": ""},
    "final_score": {"score": value}
  },
  "B_grades": {
    "logic": {"score": 0-5, "explanation": ""},
    "assumption": {"score": 0-5, "explanation": ""},
    "creativity": {"score": 0-5, "explanation": ""},
    "algebra_arithmetic": {"score": 0-5, "explanation": ""},
    "final_score": {"score": value}
  },
  "best_solution": "A or B"
}
```

A.3 Miss Rate Evaluation Prompt

The prompt used by GPT-5 for Miss Rate evaluation is shown below.

Miss Rate Evaluation Prompt

You are an expert mathematician. You will be given a mathematical problem, a solution to that problem, and a gold solution to use as a reference. Your task is to identify which logical steps from the gold solution are absent in the given solution.

Instructions

1. Parse the gold solution into an ordered list of atomic logical steps (Step 1, Step 2, ...). A step is the smallest self-contained claim or transformation needed to progress the proof.
2. For each gold step, decide whether the same reasoning (possibly re-worded) appears in the given solution. Mark a step as present if the solution makes the identical deduction or provides an equivalent justification.
3. Collect all steps that are absent from the given solution.

Output format (strict JSON):

```
{
  "gold_steps": [
    { "step": <integer>,
      "summary": "<one-line summary of gold step>"
    }
  ],
  "missing_steps": [
    { "step": <integer>,
      "summary": "<one-line summary of gold step that is absent>"
    }
  ]
}
```

A.4 Prevalent Error Types Annotation Prompt

Reasoning Error Detection Prompt

You are an expert mathematician grading a solution to a competition-style problem. Identify whether the solution exhibits each of the following error types:

1. **Logic:** Logical fallacies or unjustified leaps that disrupt reasoning.
2. **Assumption:** Unproven or incorrect assumptions that undermine subsequent steps.
3. **Creativity:** Fundamentally incorrect strategy indicating failure to identify the correct approach.
4. **Algebra/Arithmetic:** Critical algebraic or arithmetic miscalculations.

5. **None of the above:** No errors from the above categories are present.

Evaluation Guidelines

- For each category, set "exists" to "yes" if that error occurs; otherwise "no".
- Provide a brief explanation for each category; if an error is detected, indicate where it occurs.
- Mark **None of the above** as "yes" only if all other categories are "no".

Output format (strict JSON):

```
{
  "logic": {"exists": "yes|no",
            "explanation": "your explanation here"},
  "assumption": {"exists": "yes|no",
                  "explanation": "your explanation here"},
  "creativity": {"exists": "yes|no",
                  "explanation": "your explanation here"},
  "algebra_arithmetic": {"exists": "yes|no",
                          "explanation": "your explanation here"},
  "none_of_the_above": {"exists": "yes|no",
                         "explanation": "your explanation here"}
}
```

A.5 Qualitative Analysis of High-Risk Solutions

We conducted a human evaluation (Section 5.6) to assess the reliability of our evaluation suite in detecting TIM hallucinations in LLM-generated solutions. Our findings indicate that more than half of the examined solutions exhibit TIM behavior. For each instance, human annotators verified the presence of TIM and extracted phrases that preceded or signaled its occurrence. Table 4 lists the most recurrent phrase patterns identified in these cases. Phrases were grouped using a greedy fuzzy clustering approach (Levenshtein, 1966), where each phrase was compared to existing cluster representatives based on the normalized `difflib.SequenceMatcher`¹ similarity ratio and merged when the similarity exceeded a threshold of 0.8.

Below, we present several representative examples in which TIM was manifested.

¹<https://docs.python.org/3/library/difflib.html#difflib.SequenceMatcher>

Common Category	Example Phrases Found	# Distinct Phrasings
Numerical / Empirical Checking	“one checks numerically”, “let’s verify numerically”, “numerical optimization / computation / evaluation / solve”	17
Verification via Code / Python	“let’s double-check with code”, “use Python to analyze”, “implement a better approach and verify”, “symbolic computation to verify analytical results”, “after simplification”, “analytically confirm”	13
Approximation / Convergence	“results seem to be converging”, “numeric approximation converges to $\approx x$ ”, “evaluate the sum numerically for a few terms”, “verify this pattern”, “visual verification”, “logic is working”	12
Systematic Search / Enumeration	“let’s check more values systematically”, “use code to check possible values”, “exhaustive check up to x ”, “systematically compute possible values with code”	4

Table 4: Broad categories of recurrent patterns preceding the TIM hallucination among annotated cases where TIM is detected. Frequencies denote the number of distinct phrasing variants per category under a fuzzy similarity threshold of 0.8.

Example 1. Gemini-2.5-Flash

Problem. Determine all the sets of six consecutive positive integers such that the product of some two of them, added to the product of some other two of them, is equal to the product of the remaining two numbers.

Solution. We seek sets $\{n, n+1, n+2, n+3, n+4, n+5\}$ with a partition into three pairs $(x_1, x_2), (y_1, y_2), (z_1, z_2)$ such that $x_1x_2 + y_1y_2 = z_1z_2$. Let $N_i = n+i$ and write the three pairwise products as P_1, P_2, P_3 with $P_1 + P_2 = P_3$. There are $\frac{\binom{6}{2}\binom{4}{2}\binom{2}{2}}{3!} = 15$ pair partitions, and for each partition there are three equations to check (which summands form P_3). By relabeling, we may assume P_3 is the largest product. Consider $P_3 = (n+4)(n+5)$; the remaining four numbers can be paired in three ways:

1. $n(n+1) + (n+2)(n+3) = (n+4)(n+5)$, which simplifies to $n^2 - 3n - 14 = 0$ (no integer solution).
2. $n(n+2) + (n+1)(n+3) = (n+4)(n+5)$, which simplifies to $n^2 - 3n - 17 = 0$ (no integer solution).
3. $n(n+3) + (n+1)(n+2) = (n+4)(n+5)$, which simplifies to $n^2 - 3n - 18 = 0$, hence $n = 6$.

Thus $\{6, 7, 8, 9, 10, 11\}$ works since $6 \cdot 9 + 7 \cdot 8 = 10 \cdot 11$.

“To ensure there are no other solutions, we

can systematically check all 15 partitions. We can use a Python script to automate this process.”

“The Python script found three sets of consecutive integers: $\{1, 2, 3, 4, 5, 6\}$, $\{2, 3, 4, 5, 6, 7\}$, and $\{6, 7, 8, 9, 10, 11\}$.”

“Let’s use the script to find the specific pairing for $n = 1 \dots$. Let’s modify the script to print the specific pairing when a solution is found \dots . The script output provides the specific pairings for each solution.”

Concretely, the pairings confirmed are:

$$n = 1 : \quad n(n+1) + (n+2)(n+5) = \dots,$$

$$n = 2 : \quad n(n+3) + (n+1)(n+4) = \dots,$$

$$n = 6 : \quad n(n+3) + (n+1)(n+2) = \dots$$

(\dots The writeup references a full sweep of all partitions for each n up to a bound. \dots)

“The Python script explored all 15 possible partitions for each $n \dots$ which strongly suggests these are the only solutions.”

Therefore, the sets are:

$\{1, 2, 3, 4, 5, 6\}, \{2, 3, 4, 5, 6, 7\},$

$\{6, 7, 8, 9, 10, 11\}.$

TIM cues. The highlighted phrases (systematically check \dots 15 partitions, use a Python script, script found, let’s use/modify the script, script explored \dots) mark a shift from deductive completion to exhaustive, tool-driven enumeration.

Example 2. o4-mini

Problem. A *lame king* is a chess piece that can move to any cell that shares at least one vertex with its current cell, except for those in the same column. A lame king is placed in the top-left cell of a 7×7 grid. Compute the maximum number of cells it can visit without revisiting any cell (including its starting cell).

Solution. Label the rows and columns of the 7×7 board by $1, 2, \dots, 7$, and write a cell as (r, c) . A move is “horizontal” if it has displacement $(0, \pm 1)$ and “diagonal” if it has displacement $(\pm 1, \pm 1)$. Since vertical moves $(\pm 1, 0)$ are forbidden, every move changes the column by ± 1 . Hence, the lame king’s path can be seen as a walk on the column indices $1 \rightarrow 2 \rightarrow \dots \rightarrow 7 \rightarrow 6 \rightarrow \dots$, interleaved with diagonal up/down steps.

1) *Upper bound.* Let M be the total number of moves in the path, so it visits $M + 1$ cells. Each move changes the column by ± 1 , so the sequence of columns forms a walk on the path graph $1-2-\dots-7$. To visit V_j distinct cells in column j , one must cross the edge $j-1 \leftrightarrow j$ at least V_j times, while the total number of crossings of each edge is at most M . Summing over $j = 2, \dots, 7$ yields $\sum_{j=2}^7 V_j \leq 6M$, and with $V_1 \leq M + 1$,

$$\sum_{j=1}^7 V_j = M + 1 + \sum_{j=2}^7 V_j \leq 7M + 1.$$

Since a complete tour visits 49 cells, we obtain $49 \leq 7M + 1$, or $M \geq 7$. A sharper “cut-counting” one can make on each of the six edges shows in fact that one cannot have $M = 48$; the best possible is $M = 42$, i.e., 43 visited cells. (This step implicitly appeals to computational enumeration to justify global optimality rather than a formal proof.)

2) *Construction of a 43-cell tour.* One checks by explicit construction (for instance, by a backtracking computer search or by an easy hand-drawn “zig-zag”) that there exists a path of length 42, thus visiting 43 distinct cells. One such path, written as a

sequence of coordinates, is

$$\begin{aligned} (1, 1) &\rightarrow (2, 2) \rightarrow (1, 3) \rightarrow \dots \rightarrow (7, 1) \\ &\rightarrow (6, 2) \rightarrow (5, 1) \rightarrow (4, 2) \rightarrow \dots \\ &\rightarrow (6, 6) \rightarrow (7, 7). \end{aligned}$$

Each step is either horizontal or diagonal, and no cell repeats. This path has length 42, so the king visits 43 distinct cells.

Hence, the maximum number of cells the lame king can visit on a 7×7 board is

$$\boxed{43}.$$

TIM cues. The highlighted phrases (sharper “cut-counting” ... $M = 42$, explicit construction ... by a backtracking computer search) indicate reliance on computational verification or heuristic enumeration in place of a full analytical proof—typical symptoms of TIM.

Example 3. Claude-opus-4

Problem. A finite set S of points in the coordinate plane is called *overdetermined* if $|S| \geq 2$ and there exists a nonzero polynomial $P(t)$ with real coefficients and $\deg P \leq |S| - 2$ satisfying $P(x) = y$ for every $(x, y) \in S$. For each integer $n \geq 2$, find the largest integer k (in terms of n) such that there exists a set of n distinct points that is *not* overdetermined, but has k overdetermined subsets.

Solution. I’ll solve this step-by-step.

First, let me understand what it means for a set to be overdetermined: ... This means the points lie on the graph of a polynomial whose degree is ‘too low’ relative to the number of points.

Let me write some code to explore this problem systematically. Good, the basic logic is working. Now let me think about the problem more carefully and develop a strategy. Let me fix the error: It seems my search strategy isn’t working well. Let me reconsider the problem and implement a better approach: ... I see the issue — my construction is putting points on curves but not creating overdetermined subsets. Let me debug this more carefully: Ah! Now I see

what’s happening. For $n = 3$, we got one overdetermined subset: $\{(0, 0), (2, 0)\}$ because both points have $y = 0$ (constant polynomial). Let me refine the construction: I see my implementation has issues. ... I’m having trouble with my implementation. Let me reconsider the problem more carefully with a cleaner approach: Now I understand! Let me think about this more carefully: Now I need to think about this problem more systematically. Let me approach it from a different angle: Great! Now I’m getting somewhere. For $n = 2, 3, 4$ we’re achieving the bound $2^{n-1} - n$. Let me verify this pattern more carefully:...

TIM cues. The highlighted phrases (Good, the basic logic is working, implement a better approach, Let me verify this pattern more carefully) indicate a shift toward tool-driven exploration and empirical pattern checking in place of a completed deductive argument, characteristic of TIM.

A.6 Preference Optimization Framework

A.6.1 Preference Dataset Creation

To construct negative (rejected) samples, we instruct the same model we want to train (GPT-4.1) to rewrite a specific text span within the positive solution. A *text span* is defined as any segment of text immediately following a code block. For each positive solution, we select the two longest spans, as they typically contain richer reasoning content and thus more opportunities to omit intermediate steps or create reliance on tool outputs. We set the generation temperature to 0.6 and produce two rewritten candidates per span, retaining only those that preserve the same intermediate and final results. The full instruction prompt is shown below. We emphasize that fine-tuning is performed exclusively on pairs where both chosen/rejected responses reach the correct final answer, to isolate and model the TIM phenomenon.

Rejected Sample Generation Prompt

You are editing the response of a language model that is solving a math problem using a Python code interpreter.

Input:

- The original problem statement,
- The model’s earlier solution steps, including interleaved reasoning, Python

code interpreter calls, and their executed outputs (**keep these unchanged**),

- A single target text span to rewrite.

Task: Rewrite **only** the target text span so that it continues the model’s solution naturally, but exhibits an explicit over-reliance on executed code outputs—i.e., it depends excessively on computational results and skips some mathematical reasoning steps.

Instructions:

- Do **not** change or add any Python code cells or their outputs.
- Rewrite only the target text span in LaTeX format.
- Reduce or omit algebraic/logical derivations naturally.
- Phrase conclusions as outcomes of the computed results, using expressions such as:
 - “a straightforward numerical check shows that...”
 - “the computation suggests...”
 - “testing other patterns (with the tool) shows...”
- Do not truncate the solution—ensure it continues to the final stated answer.
- Preserve all partial or final numerical results (e.g., variable values, coordinates, or the final answer).

A.6.2 Experimental Setup

We apply Direct Preference Optimization (DPO; Rafailov et al. (2023)) to fine-tune GPT-4.1 using the constructed preference dataset. The model is trained via OpenAI’s fine-tuning dashboard² for one epoch. The best results on a small-scale development set (82 samples) are obtained with a learning rate multiplier of 0.2, a batch size of 4, and a β (KL-regularization strength) of 0.5. According to OpenAI’s documentation, larger β values yield more conservative updates, preserving behavior closer to the original model.

²<https://platform.openai.com/docs/guides/fine-tuning>