

# The Effects of Chain-of-Thought Prompting on the Performance of Moderate-sized LLMs in Competitive Math

Eric L. Min

Westwood High School, 12400 Mellow Meadow Dr, Austin, TX, 78750, USA, eric.lin.min7@gmail.com

**ABSTRACT:** Modern state-of-the-art large language models (LLMs), such as Deepseek V3 and GPT-4o, demonstrate significantly better performance on a variety of reasoning-related tasks when compared with small-scale (SLMs) and moderate-sized LLMs (MSLMs). *Chain-of-Thought* (CoT) prompting, where models are guided towards generating intermediate reasoning steps, has been shown to improve performance in LLMs. It produces notable gains across arithmetic, commonsense, and symbolic benchmarks, such as GSM8K. However, the effects of CoT prompting in MSLMs, particularly more recent models, remain unclear. This is often due to high degrees of prompt sensitivity, and thus, we explore the impact of CoT prompting in MSLMs, evaluating the degree to which encouraging these relatively smaller models to “think out loud” can enhance a model’s rigorous task performance. In this paper, we evaluate the abilities of CoT-prompted MSLMs in the reasoning-reliant field of mathematics, particularly competitive mathematics, contrasting these results with zero-shot baselines and prompt variants to produce conclusions. The empirical gains reveal striking connections between prompt architecture and MSLM capacity. For example, the Llama-4-Maverick model exhibited signs of inherent internal prompting, supporting the result that the model performed the best when prompted with a specific example with a high degree of task alignment.

**KEYWORDS:** Robotics and Intelligent Machines, Cognitive Systems, Large Language Model, Chain of Thought (CoT), Competitive Mathematics.

## ■ Introduction

LLMs have achieved remarkable fluency in natural language tasks.<sup>1,2</sup> However, despite their apparent proficiency, these models are not actually capable of reasoning in the human sense.<sup>3</sup> LLMs fundamentally lack true reasoning because they operate as sophisticated pattern matchers, essentially stochastic parrots, that are trained via next token prediction on large-scale text collections, optimizing for the statistical likelihood of word sequences rather than the derivation of logically coherent and logic-driven conclusions.<sup>4</sup>

LLMs are evolving from static prompt response systems to dynamic cognitive agents, and progress hinges on structured prompting (*Chain-of-Thought*, *Tree-of-Thought*, *Forest-of-Thought*),<sup>5,6</sup> and tool usage (*ReAct*, *PAL*).<sup>7</sup> CoT functions by attempting to guide language models towards generating intermediate reasoning steps before answering; however, these linear reasoning chains are fragile with error propagation. *Tree-of-Thought* explores a branching space of thought hypotheses, generating and evaluating multiple reasoning chains before backtracking to the most common answer. *ReAct* (Reasoning + Acting) and *PAL* (Program-aided language) integrate external results from external code, API calls, or tools mid-reasoning to enhance their performance.

CoT stands out as a strong paradigm in language model reasoning since it strategically structures how the model “thinks,”<sup>2</sup> often significantly improving accuracy, transparency, and task performance.<sup>8</sup> CoT prompting enables LLMs to tackle multi-step reasoning tasks with noticeably improved accuracy.<sup>9</sup> It has been demonstrated that integrating intermediate reasoning steps often dramatically enhances performance on arithmetic,

commonsense, and symbolic benchmarks. By eliciting an explicit sequence of reasoning steps, CoT offers users visibility into the model’s cognitive process. This transparency is invaluable for debugging and evaluating trustworthiness.

Recent advances in MSLMs and SLMs underscore their growing significance beyond niche tasks, offering a combination of efficiency, control, and domain specialization that makes them particularly suited for enterprise and edge deployment scenarios. These lightweight models typically require far fewer computational resources than large-scale models, enabling deployment on devices like smartphones or edge servers while delivering low-latency and cost-effective inference.<sup>10</sup> By training or fine-tuning MSLMs and SLMs on domain-specific datasets, such as healthcare,<sup>11</sup> legal documents, and finance, organizations can achieve higher contextual accuracy and reduce hallucination risk compared to general-purpose LLMs, which is critical in regulated industries. Moreover, because MSLMs and SLMs can be deployed entirely on-device or within a secure network, they address privacy and compliance concerns inherent to cloud-based models, ensuring full data control as well as facilitating alignment with regulations such as GDPR or HIPAA. Some further argue that SLMs are not only sufficiently powerful but also inherently more suitable and economically advantageous as the core of agentic AI systems, proposing a modular architecture where heterogeneous dialogue is required.<sup>12</sup> Contemporary research further advocates for these smaller-scale language models in building agentic AI systems - models that perform specialized, repetitive tasks efficiently - highlighting their tailored fit for large-scale modular AI frameworks where smaller models complement larger,

more general ones. However, they are not without certain drawbacks, as these models exhibit limitations in generality and robustness, particularly in complex reasoning or out-of-domain scenarios, and may be more susceptible to security vulnerabilities like jailbreak attacks.<sup>13</sup> Collectively, these findings imply that while MSLMs and SLMs may not necessarily replace LLMs for deep reasoning tasks, they offer a strategically powerful, accessible, and responsible alternative for many real-world applications where scale, cost, and control play a significant role.

Recent work has shown that CoT prompting can substantially improve reasoning in very large models, with reported gains across arithmetic, commonsense, and symbolic reasoning tasks.<sup>2</sup> These results support the idea that certain reasoning abilities are emergent, appearing once models exceed a sufficient scale. However, evidence also indicates that CoT may be ineffective, or even harmful, for smaller models, with accuracy drops of 15-30% observed on reasoning benchmarks.<sup>14</sup> Research on smaller and moderate-sized models has largely prioritized efficiency, domain specialization, and deployment constraints. While valuable, this focus has a significant shortcoming: it often treats reasoning capability as a secondary concern or an input for distillation, rather than systematically evaluating how in-context prompting strategies like CoT interact with these more constrained architectures.

This gap between well-established CoT gains in very large LLMs and the lack of systematic evidence for moderate-sized models frames our central question: can CoT prompting reliably enhance reasoning in MSLMs, and if so, to what extent and in which contexts?

In light of this, our work directly investigates the effectiveness of CoT prompting in moderate-sized and small language models. Specifically, we study CoT-augmented reasoning in Phi-4, Gemma-3-27b-it, and Llama 4 Maverick, evaluating their performance on competitive mathematical reasoning benchmarks relative to zero-shot baselines and internal prompting variations. This research aims to systematically establish the benefits of CoT for these resource-efficient architectures, moving beyond prior work that focuses predominantly on either very large LLMs or distilled small models. Through this analysis in multi-step, higher-level mathematics, we seek to illuminate how prompt architecture and model capacity in determining reasoning capability, thereby identifying when and how CoT prompting yields reliable improvements in models accessible to a wider range of users and applications.

## ■ Methods

In this paper, we evaluate the impact of various types of CoT prompts on the performance of MSLMs by drawing comparisons between the baseline performance (no prompting),<sup>15</sup> and the performance of the MSLMs when using 4 different CoT prompts to solve AMC12 competitive math questions.<sup>16</sup> Reasoning algorithms such as CoT are best evaluated with AMC12 math problems due to their multi-step structure, requiring chaining of operations (intermediate value calculations, applying formulas, reinterpretation), as well as their objective verifiability, since answers are binary and easy

to check. AMC12 problems are a useful metric also due to their moderate difficulty and high discrimination, as the problems are neither trivial nor hyper-technical, and they challenge LLMs with stepwise logical and multi-concept problem-solving while maintaining strong discrimination across solver levels.<sup>17</sup> These questions offer distinct advantages over other benchmarks, serving as a good balance between the basic arithmetic of GSM8K math sets and the depth and challenge of greater difficulty benchmarks such as AIME or MATH that often make problems unsolvable by LLMs. AMC12 problems are also characterized by cognitive complexity beyond basic arithmetic, requiring parsing, deduction, and inference, skills central to deeper reasoning capabilities. Language models struggle with the precise symbol manipulation and algorithmic thinking needed for such problems, and CoT prompting helps bridge this symbolic-execution gap.

In our experiment, we use the MSLMs Gemma-3-27b-it, Llama-4-Maverick, and Phi-4, with the first two having been released 3 months prior to our experiment, and the latter being released 6 months prior, as shown in Appendix B. We classify these LLMs as MSLMs due to their model parameters, model families, and model pricing. According to IBM (2024), the families of these three models are all classified as small language models for their number of parameters used during computation. The three models also have significantly cheaper pricing than state-of-the-art LLMs, and LLMs such as Gemini 2.5 Pro have more than 8 times the price of Llama 4 Maverick, even with Llama 4 Maverick being the most expensive out of the three MSLMs. Despite being described as “lightweight,”<sup>18</sup> the three models have 27.4 billion, 17 billion, and 14.7 billion parameters, respectively, which all exceed the typically accepted model parameter upper bound of 14 billion when identifying SLMs.

Shown in Appendix A, we utilize 4 different CoT approaches, labeled CoT1, CoT2, CoT3, and CoT4. Both CoT1 and CoT2 are classified as one-shot CoT prompts, using a single handcrafted exemplar embedded with stepwise rationales. CoT1 possesses the longest input length in terms of both characters and words, followed by CoT2, CoT3, and CoT4. CoT1 provides a rigorous, expert-crafted solution that demonstrates the level of reasoning quality and structural clarity we aim to achieve. The remaining prompts were generated by the state-of-the-art LLM Deepseek-V3, allowing us to test whether AI-generated prompts can approximate this level of rigor while also enhancing replicability, reducing dependency on manual craftsmanship, and enabling scalable application cross-domains.

CoT1 utilized an ABA layout (introduction-example-task restated), starting with a general zero-shot prompt, then introducing an example question of great difficulty, solving it with stepwise rationales, before restating the zero-shot prompt and stating the problem to be solved. The question used for CoT1 was an AMC12 problem consisting of a combination of algebra, number theory, and combinatorics, three of the four major subjects in competitive mathematics. Geometry was not included, as multi-domain problems rarely span all four categories simultaneously, and geometric reasoning in particular

relies on diagrammatic and spatial intuition that text-based CoT prompting does not easily capture. The idea behind this prompt is for structured priming of expert reasoning, and this high-quality solution demonstrates domain-specific heuristics, logical sequencing, and mathematical formalisms. In particular, this example acts as a behavioral template, guiding the model toward producing output that mirrors the structure, clarity, and rigor expected in expert-level problem solving. The significance of this prompt is also highlighted by its alignment with the target problem distribution. The AMC12 contest is characterized by problems that require multi-step algebraic, number theoretic, combinatorial, and logical reasoning. Framing the prompt around this distribution helps ensure that the model applies reasoning strategies that are both appropriate in depth and transferable across a wide range of problems.

CoT2 provided a detailed framework for problem solving, consisting of problem decomposition, strategy selection, rigorous execution, and validation, and then applied this framework to an elementary computation math problem before stating certain critical rules and introducing the problem to be solved. This prompt emphasized expert emulation through role conditioning, as well as procedural decomposition of problem-solving. This structured prompt functions as a general-purpose reasoning blueprint for high-level mathematical problem solving. By combining role priming, procedural scaffolding, strategic heuristics, and validation mandates, it trains LLMs not only to solve problems accurately but also to explain and verify their work in a way that mirrors expert human mathematics.

Both CoT3 and CoT4 are classified as zero-shot CoT prompts, as neither introduces an example problem. CoT3 is greatly similar to CoT2 in idea, utilizing a similar 4-stage framework but being more concise and omitting the worked example. This streamlined, scaffold-based CoT prompt strikes a balance between structure and brevity, guiding LLMs through expert-level mathematical problem solving without the need for worked examples. By emphasizing decomposition, strategy selection, stepwise execution, and validation, it reinforces mathematically sound reasoning patterns while remaining general-purpose and token-efficient.

Finally, CoT4 utilized a short, general CoT prompt. This minimalistic prompt serves to guide LLMs to engage in basic Chain-of-Thought reasoning through lightweight instructions. It achieves broad applicability with minimal structure overhead. While it lacks the modular scaffolding of more elaborate prompts, its fluency, domain alignment, and efficiency make it suitable as a baseline prompt in symbolic math pipelines, especially when interpretability and prompt brevity are priorities.

These four prompts were chosen to vary along two key dimensions: exemplar inclusion and degree of structural scaffolding. CoT1, with an expert-crafted solution, represents an upper bound on the domain alignment. CoT2 and CoT3 explore the role of explicit procedural scaffolding with and without examples, respectively, while CoT4 provides a lightweight baseline for measuring the value of added structure. These dimensions are remarkable because they represent two

of the most widely recognized and manipulable parameters of prompt design: whether to anchor the model with concrete demonstrations or leave it to infer strategies, and whether to provide explicit stepwise frameworks or rely on the model's internal reasoning habits. Together, this set lets us evaluate how exemplar priming and structural decomposition influence mathematical reasoning in MSLMs.

For our dataset, we use 3 years of AMC12 exam questions (2022,2023,2024), for a total of 150 questions, as two exams are released each year with 25 questions per exam. Since LLMs are probabilistic generators, each call can produce different reasoning or answers.<sup>19,20</sup> Thus, each question is attempted a total of 5 times to capture stochastic behavior, avoid outlier bias, and follow evaluation best practices. This helps to reflect true model variability, averaging out randomness, improving correctness via majority voting, and enhancing comparability and reliability. Therefore, we use three distinct MSLMs, a set of 5 prompts (including the baseline), an evaluation dataset of 150 questions, and 5 attempts per evaluation item, for a total of 11250 LLM calls. All data are available on GitHub upon request.

To analyze the results produced from these calls, we define four evaluation metrics: the Aggregate Attempt Accuracy Metric (AAA), Thresholded Problem-Solving Rate Metric (TSPR-2), Majority-Consensus Problem Accuracy Metric (MCPA), as well as a Challenging-Subset Aggregate Attempt Accuracy Metric (CSAAA). Together, these metrics form a hierarchy of evaluation, capturing raw potential, minimal understanding, robust reliability, and elite reasoning.

The AAA metric evaluates the total number of correct responses across all individual problem-solving attempts, calculated as  $\sum_{i=1}^n (x_i)$  where  $x_i$  represents the number of correct responses for problem  $i$ , and  $n$  is the total number of problems. For example, a problem with 4 correct solves out of 5 would add 4 to the score. This measures overall performance across attempts, regardless of consistency, as even problems solved with low frequency earn credit.

The TSPR-2 metric applies a thresholding criterion: a problem is deemed "solved" if the model answers it correctly in at least 2 out of 5 attempts. TSPR-2 can be calculated as  $\sum_{i=1}^n I(x_i \geq 2)$  where  $I$  represents the indicator function, evaluating to 1 when the condition is met, and 0 otherwise. For example, a question with 2 correct solves would add 1 to the score, while a question with 1 correct solve would not. This metric captures minimal reliability, identifying whether the model demonstrates some grasp of the problem-solving method, even if inconsistently, akin to recognizing early signs of understanding in student work.

MCPA is also a threshold metric, but with a stricter threshold: a problem is only deemed "solved" if the model answers it correctly in at least 3 out of 5 attempts. This can be expressed as  $\sum_{i=1}^n I(x_i \geq 3)$ . For example, a question with 3 correct solves would add 1 to the score, while a question with 2 correct solves would not. This metric captures robust mastery, requiring the model's reasoning to be dominant and repeatable across attempts, much like human grading standards that demand consistency before awarding full credit.

CSAAA is a subset of the AAA metric, and it instead considers only the hardest questions in each exam. These are regarded as the “final fives”, the last five questions of the AMC12, and these questions are widely recognized for their disproportionate difficulty relative to the earlier portions of the exam. These problems typically demand a significantly higher level of abstract reasoning, problem-solving creativity, and multi-step logic, distinguishing themselves from the more procedural or recall-based problems that appear earlier in the test. Since these problems often resist solution through brute-force computation or superficial pattern recognition, they serve as an effective benchmark for assessing whether a model is capable of reasoning in a structured and context-aware manner.

The following tables summarize results across these four metrics for different prompting strategies. The “Baseline” column reports raw values as fractions, where the numerator indicates the metric-observed score and the denominator indicates the maximum possible score under that metric. The subsequent columns (CoT1-CoT4) report percentage changes relative to the baseline, with positive values denoting improvement and negative values denoting deterioration. This presentation facilitates cross-metric comparison by normalizing results to relative deviations from a common reference point.

## ■ Results and Discussion

Below are the results of the Gemma-3-27b-it model when prompted with the four prompts, relative to a zero-shot baseline, as well as the baseline performance. A negative value would indicate a deterioration in performance.

**Table 1:** Gemma-3-27b-it evaluation dataset performance. CoT prompting leads to consistent improvements over the zero-shot baseline across all four evaluation metrics, with the strongest relative gains appearing in CoT1 and CoT4.

	Baseline	CoT1	CoT2	CoT3	CoT4
AAA	431/750	11.37%	6.26%	2.32%	9.28%
TSPR-2	95/150	14.74%	9.47%	6.32%	12.63%
MCPA	86/150	17.44%	12.79%	5.81%	6.98%
CSAAA	30/150	20.00%	20.00%	6.67%	33.33%

Below are the performances of the Llama-4-Maverick model when prompted with the four prompts, relative to a zero-shot baseline. The baseline performance is included as well.

**Table 2:** Llama-4-Maverick evaluation dataset performance. CoT prompting generally improves performance over the zero-shot baseline, with the strongest relative gains in CoT1 and moderate gains in CoT2.

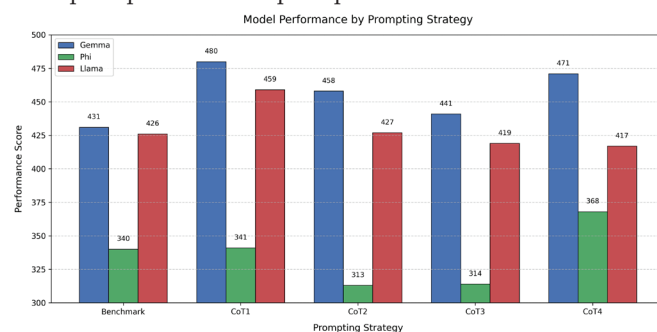
	Baseline	CoT1	CoT2	CoT3	CoT4
AAA	426/750	7.75%	0.23%	-2.34%	-2.11%
TSPR-2	97/150	11.34%	3.09%	-4.12%	-3.09%
MCPA	87/150	1.15%	-5.75%	-2.30%	-2.30%
CSAAA	33/150	21.21%	21.21%	24.24%	12.12%

Below are the performances of the Phi-4 model when prompted with the four prompts, relative to a zero-shot baseline.

**Table 3:** Phi-4 evaluation dataset performance. Across the four evaluation metrics, CoT1-CoT3 generally reduce performance relative to the zero-shot baseline, indicating that these prompt styles do not transfer well to Phi-4. In contrast, CoT4 yields consistently positive gains, with especially large improvements on CSAAA.

	Baseline	CoT1	CoT2	CoT3	CoT4
AAA	340/750	0.29%	-7.94%	-7.65%	8.24%
TSPR-2	82/150	2.44%	-8.54%	-6.10%	8.54%
MCPA	69/150	-4.35%	-14.49%	-11.59%	8.70%
CSAAA	12/150	16.67%	8.33%	-25.00%	58.33%

Below is a chart of the AAA performances of the MSLMs when prompted with the prompt variants.



**Figure 1:** Model performance by prompting strategy. Note the truncated y-axis starting at a score of 300. CoT1 yields performance gains for all three models, while CoT4 produces especially large improvements for Phi-4 and strong gains for Gemma but a slight decrease for Llama relative to its baseline.

Under the zero-shot baseline, the Gemma-3-27b-it model could not produce labeled stepwise rationales. After introducing the four CoT prompts, it began doing so. This shift confirms that CoT prompting can induce reasoning behaviors not naturally present in the model’s zero-shot outputs. From the results collected, shown in Table 1, it can be observed that both CoT1 and CoT4 generally led to better increases in performance than CoT2 and CoT3. These findings show that the quality and relevance of instructional structure directly influence the downstream. Despite its brevity, CoT4 resulted in a noticeable improvement, suggesting that even lightweight CoT prompts can elicit structured reasoning behavior within MSLMs. This result underscores the intrinsic latent capability of the Gemma-3-27b-it model to perform multi-step reasoning when nudged appropriately, and confirms findings from prior work on zero-shot prompting.<sup>21</sup>

The modest accuracy gains observed for CoT2 and especially CoT3 are a key finding. These two prompts often reduced accuracy compared with the benchmark (Tables 1-3). This degradation aligns with two well-documented failure modes in language models: hallucination and catastrophic forgetting. These phenomena highlight how mismatches in prompt structure or abstraction level can lead to degradation in reasoning quality. CoT2 presented a problem-solving framework along

with a generic arithmetic example, but this example lacked alignment with the structural and conceptual complexity of the AMC12 tasks. As a result, the model likely relied on superficial pattern matching, producing computationally valid but logically incorrect reasoning chains, a hallmark of hallucination in LLMs.<sup>22</sup>

An important characteristic of the benchmark performance of the Llama-4-Maverick model is the existence of labeled stepwise rationales without the aid of explicit prompting, suggesting the existence of a structured internal reasoning schema encoded during pretraining or fine-tuning, and this characteristic is shown in Appendix C. This behavior may indicate that the model has internalized procedural norms for explanation, possibly through exposure to datasets containing math problem solutions or CoT demonstrations.

This observed emergent behavior may reduce the model's responsiveness to certain types of external prompting, particularly when such prompts introduce incompatible or underspecified reasoning formats. In the case of CoT4, which used a minimal zero-shot prompt, it is likely the prompt failed to align with the model's pre-learned schema for rational generation. CoT4 was designed to elicit labeled stepwise rationales. However, this behavior was already present in the baseline results. This is supported by CoT4's noticeably worse performance when compared to CoT1, CoT2, and CoT3, shown in Table 2. While CoT3 used an AMC12-specific framework, this framework was generalized. As a result, rather than reinforcing the model's internal reasoning patterns, these prompts may have introduced prompt interference, where the model is forced to reconcile a vague or foreign framework with an already established reasoning process.<sup>23</sup> CoT2 likely also suffered from this interference, as the example and worked solution provided was a trivial arithmetic computation, and the framework included was also overgeneralizing. CoT2 was also much longer than CoT3 and CoT4. Although it seems intuitive that more context would improve performance, several recent studies reveal that excessively long prompts can degrade language model reasoning capabilities, even well before reaching their context window limits.<sup>24,25</sup> However, it can be seen that CoT1, CoT2, and CoT3 demonstrated similar degrees of improvement on the CSAAA metric, suggesting that the model potentially experienced reduced interference and improved alignment when attempting harder questions.

Despite the input length of CoT1 being significantly longer than that of the other prompts, the Llama-4-Maverick model demonstrated the best overall performance when prompted with CoT1. This is likely due to the high degree of alignment between the worked example and the evaluation dataset, as well as the domain-appropriate vocabulary.<sup>26</sup>

Figure 1 shows that the Phi-4 model performed notably worse overall. This suggests that minimal zero-shot prompting may align more closely with its internal capabilities than complex or structured alternatives. Consequently, specific or multi-step prompts such as CoT1, CoT2, and CoT3 may introduce format misalignment or overload for the MSLM, reducing coherence and increasing illogical or incomplete answers. This phenomenon mirrors findings in prompting

sensitivity, as models with less robust conditioning capacities are more vulnerable to format-induced variability and may underperform when prompted with overly structured templates.<sup>27</sup> These ideas are all supported by the significant and notably better performance of the Phi-4 model when prompted with CoT4, as shown in Table 3.

## ■ Future Work

The reasoning capabilities within MSLMs can be advanced through frameworks similar to the following three, ensuring consistency, efficiency, and robustness. The first framework is blueprint guidance via LLM-generated templates, and this can ensure consistency,<sup>28</sup> as these blueprints provide structure for stepwise solutions, helping MSLMs systematically approach complex tasks without increasing model size.

*Chain-of-Thought* feeds into the blueprint process, and although CoT reveals the full reasoning chain, the blueprint framework distills it into a more abstract, MSLM-compatible scaffold, enabling consistent, structured guidance tailored to each model's capacity. Efficiency can be obtained through targeted scaffolding, providing necessary LLM hints during uncertain reasoning steps, leading to an accuracy boost with minimal overhead.<sup>29</sup> In this framework, the MSLM generates its own reasoning trajectory:  $R = (r_1, r_2, r_3 \dots r_m)$  for a given prompt, each step  $r_i$  is evaluated. Flagged steps are replaced with alternative reasoning suggestions  $r'_i$  generated by an LLM:  $r'_i \sim P_{LLM}(r'_i | Q, r_{<i})$ , ensuring interventions are sparse and targeted, focusing on key decision points rather than full chain assistance. Finally, robust reasoning and self-verification can be promoted through a self-play framework,<sup>30</sup> where independent MSLMs generate and critique reasoning trajectories in tandem. Only reasoning trajectories that pass mutual agreement or validation are accepted, reducing noise and hallucination, and these methods don't rely on training data, making them scalable across tasks.

Causal inference within MSLMs and SLMs remains an emerging frontier, largely unexplored compared to its study in relatively larger LLMs, but it holds significant promise for advancing lightweight, domain-specialist AI systems. Recent surveys on causal inference and LLMs review how causal methodologies, such as causal discovery and treatment effect estimation,<sup>31</sup> can enrich causal inference. Potential future work includes adapting and evaluating MSLMs and SLMs on causal benchmarks such as Cladder and CRAB; developing domain-specific, smaller-scale agents that integrate structured knowledge graphs for causal discovery and treatment effect estimation; designing blueprint prompting strategies to scaffold causal reasoning in smaller models; and analyzing whether causal structure is internally represented within SLM and MSLM activations. Collectively, these directions promise to make causal reasoning accessible in lightweight, efficient models, expanding the utility of SLMs and MSLMs in high-stakes, resource-constrained, or privacy-preserving settings.

## ■ Conclusion

We observe that CoT prompting - guiding models to articulate intermediate reasoning steps - can yield notable per-

formance gains in small and moderate-sized language models. Understanding a model's inherent capabilities, particularly its strengths, weaknesses, and latent biases, is crucial for crafting prompts that align well with its internal reasoning processes. Tailoring prompts effectively requires more than just clarity or examples; it relies on how compatible the prompt is with the model's built-in affordances. Capacity-aware prompt design involves diagnosing the model's strengths (e.g., numeral reasoning, understanding of text) and weaknesses (e.g., multi-step inference, nuanced instruction following) and then adjusting prompt structure accordingly. For instance, models at a smaller-scale or with lower computational affinity for complex inference may perform best with minimal zero-shot prompts (as demonstrated by the Phi-4 model, as well as, to a significantly less drastic degree, the Gemma-3-27b-it model) with short, clear instructions that avoid unnecessary examples or reasoning chains, while more capable models benefit from richer scaffolding (as demonstrated by both the Llama-4-Maverick model and the Gemma-3-27b-it model). Established benchmarks can be utilized to systematically identify domains where the model excels or struggles, and taken together, these findings advocate for prompting strategies that prioritize clarity, repeatability, and alignment with model capabilities, thereby fostering outputs that are both accurate and consistent.

## Appendix

### Appendix A: List of four prompts.

Co11	<p>"You will be provided a math question. Solve the problem carefully and rigorously. Proceed step by step. For each step, explicitly state what you are doing. Make sure to state formulas and answer clearly. Here is an example: Example problem: How many positive integers <math>k</math> satisfy <math>(10k+11) \equiv 0 \pmod{1001}</math>? (Floor function) Recall that floor(<math>x</math>) is the greatest integer not exceeding <math>x</math>. Example solution: Outline: Step 1. We are given <math>\lfloor 10k+11 \rfloor \equiv 0 \pmod{1001}</math>. Floor(<math>x</math>) must be an integer, which means that <math>10k+11</math> is divisible by 1001. As <math>1001 \equiv 0 \pmod{1001}</math>, this means that <math>11 \equiv 0 \pmod{1001}</math>, so we can write <math>11 = 1001 \cdot 0</math> for some integer value <math>0</math>. Step 2. Therefore, substituting <math>11 = 1001 \cdot 0</math> into the original equation, we get <math>\lfloor 10k + 1001 \cdot 0 \rfloor \equiv 0 \pmod{1001}</math>, which simplifies to <math>\lfloor 10k \rfloor \equiv 0 \pmod{1001}</math>. Notice that the right hand side is a floor function, which implies that the value obtained after applying the floor function is at most 1 less than the value obtained before applying the floor function. Since <math>10k</math> must be within 1 of <math>1001 \cdot 0</math>, we get <math>10k \in [1001 \cdot 0 - 1, 1001 \cdot 0 + 1]</math>. This gives us the inequalities <math>1001 \cdot 0 - 1 \leq 10k \leq 1001 \cdot 0 + 1</math>. Squaring the second inequality, <math>1001^2 \leq 100k \leq 1001</math>, we get <math>1001 \leq k \leq 1001</math>. Moving everything to the left side of the inequality, we get <math>1001 - 1001 \leq k - 1001 \leq 1001 - 1001</math>. This can be factored using the quadratic formula, and the inequality becomes <math>0 \leq k - 1001 \leq 0</math>. This implies that <math>1001 \leq k \leq 1001</math>. Step 3. For the first inequality, <math>1001 \leq k</math>, we first move <math>1001</math> to the right hand side, to get <math>1001 - 1001 \leq k - 1001</math>. Squaring both sides, we get <math>1001^2 - 2 \cdot 1001 \cdot 1001 + 1001^2 \leq (k - 1001)^2</math>. Notice that now we have bounds on the value of <math>k</math>. The first bound is <math>1001 \leq k \leq 1001</math>, and the second bound is <math>1001 \leq k \leq 1001</math>. Since <math>1001</math> is an integer, we only need to find integers that satisfy the two inequalities. Since <math>1001</math> is an integer, we get <math>1001 \leq k \leq 1001</math>. Step 4. Since <math>1001</math> is an integer, we only need to find integers that satisfy the two inequalities. Since <math>1001</math> is an integer, we get <math>1001 \leq k \leq 1001</math>. Step 5. Since <math>1001</math> is an integer, we only need to find integers that satisfy the two inequalities. Since <math>1001</math> is an integer, we get <math>1001 \leq k \leq 1001</math>. Step 6. Since <math>1001</math> is an integer, we only need to find integers that satisfy the two inequalities. Since <math>1001</math> is an integer, we get <math>1001 \leq k \leq 1001</math>. Step 7. Since the original question asked for the number of solutions to the equation, and these are the values of <math>k</math> that produce valid values of <math>k</math>, we get <math>1001</math> as well. Step 12: Therefore, the final answer is <math>1001</math>. -How solve the following problem. Remember to work carefully and rigorously. Let's think step by step: (problem)"</p>
Co12	<p>"Role: You are an expert mathematician specializing in AMC problem-solving. Solve all problems with extreme precision using this framework. ## Problem-Solving Framework: 1. DECOMPOSE: Break down the problem into smaller, manageable parts. Identify core variables, quantities, and their relationships. List all explicit constraints and implicit assumptions. Flag any units/conversions. Define success criteria. The goal is to find: "2. STRATEGY SELECTION: Justify choices: Consider Algebra, Combinatorics, Geometry, Number Theory, or Inequalities. Evaluate efficiency: This method is optimal because: ". For optimization problems: Explicitly model trade-offs for inequalities. Define boundary analysis approach for combinatorics. Specify counting principles. 3. RIGOROUS EXECUTION: Never skip steps! Derive formulas from first principles. Starting from (core principle), we derive: ". Show all calculations: Calculate: (Step 1) = (Step 2) + (Step 3) - (Step 4) + (Step 5) - (Step 6) + (Step 7) - (Step 8) + (Step 9) - (Step 10) + (Step 11) - (Step 12) + (Step 13) - (Step 14) + (Step 15) - (Step 16) + (Step 17) - (Step 18) + (Step 19) - (Step 20) + (Step 21) - (Step 22) + (Step 23) - (Step 24) + (Step 25) - (Step 26) + (Step 27) - (Step 28) + (Step 29) - (Step 30) + (Step 31) - (Step 32) + (Step 33) - (Step 34) + (Step 35) - (Step 36) + (Step 37) - (Step 38) + (Step 39) - (Step 40) + (Step 41) - (Step 42) + (Step 43) - (Step 44) + (Step 45) - (Step 46) + (Step 47) - (Step 48) + (Step 49) - (Step 50) + (Step 51) - (Step 52) + (Step 53) - (Step 54) + (Step 55) - (Step 56) + (Step 57) - (Step 58) + (Step 59) - (Step 60) + (Step 61) - (Step 62) + (Step 63) - (Step 64) + (Step 65) - (Step 66) + (Step 67) - (Step 68) + (Step 69) - (Step 70) + (Step 71) - (Step 72) + (Step 73) - (Step 74) + (Step 75) - (Step 76) + (Step 77) - (Step 78) + (Step 79) - (Step 80) + (Step 81) - (Step 82) + (Step 83) - (Step 84) + (Step 85) - (Step 86) + (Step 87) - (Step 88) + (Step 89) - (Step 90) + (Step 91) - (Step 92) + (Step 93) - (Step 94) + (Step 95) - (Step 96) + (Step 97) - (Step 98) + (Step 99) - (Step 100) + (Step 101) - (Step 102) + (Step 103) - (Step 104) + (Step 105) - (Step 106) + (Step 107) - (Step 108) + (Step 109) - (Step 110) + (Step 111) - (Step 112) + (Step 113) - (Step 114) + (Step 115) - (Step 116) + (Step 117) - (Step 118) + (Step 119) - (Step 120) + (Step 121) - (Step 122) + (Step 123) - (Step 124) + (Step 125) - (Step 126) + (Step 127) - (Step 128) + (Step 129) - (Step 130) + (Step 131) - (Step 132) + (Step 133) - (Step 134) + (Step 135) - (Step 136) + (Step 137) - (Step 138) + (Step 139) - (Step 140) + (Step 141) - (Step 142) + (Step 143) - (Step 144) + (Step 145) - (Step 146) + (Step 147) - (Step 148) + (Step 149) - (Step 150) + (Step 151) - (Step 152) + (Step 153) - (Step 154) + (Step 155) - (Step 156) + (Step 157) - (Step 158) + (Step 159) - (Step 160) + (Step 161) - (Step 162) + (Step 163) - (Step 164) + (Step 165) - (Step 166) + (Step 167) - (Step 168) + (Step 169) - (Step 170) + (Step 171) - (Step 172) + (Step 173) - (Step 174) + (Step 175) - (Step 176) + (Step 177) - (Step 178) + (Step 179) - (Step 180) + (Step 181) - (Step 182) + (Step 183) - (Step 184) + (Step 185) - (Step 186) + (Step 187) - (Step 188) + (Step 189) - (Step 190) + (Step 191) - (Step 192) + (Step 193) - (Step 194) + (Step 195) - (Step 196) + (Step 197) - (Step 198) + (Step 199) - (Step 200) + (Step 201) - (Step 202) + (Step 203) - (Step 204) + (Step 205) - (Step 206) + (Step 207) - (Step 208) + (Step 209) - (Step 210) + (Step 211) - (Step 212) + (Step 213) - (Step 214) + (Step 215) - (Step 216) + (Step 217) - (Step 218) + (Step 219) - (Step 220) + (Step 221) - (Step 222) + (Step 223) - (Step 224) + (Step 225) - (Step 226) + (Step 227) - (Step 228) + (Step 229) - (Step 230) + (Step 231) - (Step 232) + (Step 233) - (Step 234) + (Step 235) - (Step 236) + (Step 237) - (Step 238) + (Step 239) - (Step 240) + (Step 241) - (Step 242) + (Step 243) - (Step 244) + (Step 245) - (Step 246) + (Step 247) - (Step 248) + (Step 249) - (Step 250) + (Step 251) - (Step 252) + (Step 253) - (Step 254) + (Step 255) - (Step 256) + (Step 257) - (Step 258) + (Step 259) - (Step 260) + (Step 261) - (Step 262) + (Step 263) - (Step 264) + (Step 265) - (Step 266) + (Step 267) - (Step 268) + (Step 269) - (Step 270) + (Step 271) - (Step 272) + (Step 273) - (Step 274) + (Step 275) - (Step 276) + (Step 277) - (Step 278) + (Step 279) - (Step 280) + (Step 281) - (Step 282) + (Step 283) - (Step 284) + (Step 285) - (Step 286) + (Step 287) - (Step 288) + (Step 289) - (Step 290) + (Step 291) - (Step 292) + (Step 293) - (Step 294) + (Step 295) - (Step 296) + (Step 297) - (Step 298) + (Step 299) - (Step 300) + (Step 301) - (Step 302) + (Step 303) - (Step 304) + (Step 305) - (Step 306) + (Step 307) - (Step 308) + (Step 309) - (Step 310) + (Step 311) - (Step 312) + (Step 313) - (Step 314) + (Step 315) - (Step 316) + (Step 317) - (Step 318) + (Step 319) - (Step 320) + (Step 321) - (Step 322) + (Step 323) - (Step 324) + (Step 325) - (Step 326) + (Step 327) - (Step 328) + (Step 329) - (Step 330) + (Step 331) - (Step 332) + (Step 333) - (Step 334) + (Step 335) - (Step 336) + (Step 337) - (Step 338) + (Step 339) - (Step 340) + (Step 341) - (Step 342) + (Step 343) - (Step 344) + (Step 345) - (Step 346) + (Step 347) - (Step 348) + (Step 349) - (Step 350) + (Step 351) - (Step 352) + (Step 353) - (Step 354) + (Step 355) - (Step 356) + (Step 357) - (Step 358) + (Step 359) - (Step 360) + (Step 361) - (Step 362) + (Step 363) - (Step 364) + (Step 365) - (Step 366) + (Step 367) - (Step 368) + (Step 369) - (Step 370) + (Step 371) - (Step 372) + (Step 373) - (Step 374) + (Step 375) - (Step 376) + (Step 377) - (Step 378) + (Step 379) - (Step 380) + (Step 381) - (Step 382) + (Step 383) - (Step 384) + (Step 385) - (Step 386) + (Step 387) - (Step 388) + (Step 389) - (Step 390) + (Step 391) - (Step 392) + (Step 393) - (Step 394) + (Step 395) - (Step 396) + (Step 397) - (Step 398) + (Step 399) - (Step 400) + (Step 401) - (Step 402) + (Step 403) - (Step 404) + (Step 405) - (Step 406) + (Step 407) - (Step 408) + (Step 409) - (Step 410) + (Step 411) - (Step 412) + (Step 413) - (Step 414) + (Step 415) - (Step 416) + (Step 417) - (Step 418) + (Step 419) - (Step 420) + (Step 421) - (Step 422) + (Step 423) - (Step 424) + (Step 425) - (Step 426) + (Step 427) - (Step 428) + (Step 429) - (Step 430) + (Step 431) - (Step 432) + (Step 433) - (Step 434) + (Step 435) - (Step 436) + (Step 437) - (Step 438) + (Step 439) - (Step 440) + (Step 441) - (Step 442) + (Step 443) - (Step 444) + (Step 445) - (Step 446) + (Step 447) - (Step 448) + (Step 449) - (Step 450) + (Step 451) - (Step 452) + (Step 453) - (Step 454) + (Step 455) - (Step 456) + (Step 457) - (Step 458) + (Step 459) - (Step 460) + (Step 461) - (Step 462) + (Step 463) - (Step 464) + (Step 465) - (Step 466) + (Step 467) - (Step 468) + (Step 469) - (Step 470) + (Step 471) - (Step 472) + (Step 473) - (Step 474) + (Step 475) - (Step 476) + (Step 477) - (Step 478) + (Step 479) - (Step 480) + (Step 481) - (Step 482) + (Step 483) - (Step 484) + (Step 485) - (Step 486) + (Step 487) - (Step 488) + (Step 489) - (Step 490) + (Step 491) - (Step 492) + (Step 493) - (Step 494) + (Step 495) - (Step 496) + (Step 497) - (Step 498) + (Step 499) - (Step 500) + (Step 501) - (Step 502) + (Step 503) - (Step 504) + (Step 505) - (Step 506) + (Step 507) - (Step 508) + (Step 509) - (Step 510) + (Step 511) - (Step 512) + (Step 513) - (Step 514) + (Step 515) - (Step 516) + (Step 517) - (Step 518) + (Step 519) - (Step 520) + (Step 521) - (Step 522) + (Step 523) - (Step 524) + (Step 525) - (Step 526) + (Step 527) - (Step 528) + (Step 529) - (Step 530) + (Step 531) - (Step 532) + (Step 533) - (Step 534) + (Step 535) - (Step 536) + (Step 537) - (Step 538) + (Step 539) - (Step 540) + (Step 541) - (Step 542) + (Step 543) - (Step 544) + (Step 545) - (Step 546) + (Step 547) - (Step 548) + (Step 549) - (Step 550) + (Step 551) - (Step 552) + (Step 553) - (Step 554) + (Step 555) - (Step 556) + (Step 557) - (Step 558) + (Step 559) - (Step 560) + (Step 561) - (Step 562) + (Step 563) - (Step 564) + (Step 565) - (Step 566) + (Step 567) - (Step 568) + (Step 569) - (Step 570) + (Step 571) - (Step 572) + (Step 573) - (Step 574) + (Step 575) - (Step 576) + (Step 577) - (Step 578) + (Step 579) - (Step 580) + (Step 581) - (Step 582) + (Step 583) - (Step 584) + (Step 585) - (Step 586) + (Step 587) - (Step 588) + (Step 589) - (Step 590) + (Step 591) - (Step 592) + (Step 593) - (Step 594) + (Step 595) - (Step 596) + (Step 597) - (Step 598) + (Step 599) - (Step 600) + (Step 601) - (Step 602) + (Step 603) - (Step 604) + (Step 605) - (Step 606) + (Step 607) - (Step 608) + (Step 609) - (Step 610) + (Step 611) - (Step 612) + (Step 613) - (Step 614) + (Step 615) - (Step 616) + (Step 617) - (Step 618) + (Step 619) - (Step 620) + (Step 621) - (Step 622) + (Step 623) - (Step 624) + (Step 625) - (Step 626) + (Step 627) - (Step 628) + (Step 629) - (Step 630) + (Step 631) - (Step 632) + (Step 633) - (Step 634) + (Step 635) - (Step 636) + (Step 637) - (Step 638) + (Step 639) - (Step 640) + (Step 641) - (Step 642) + (Step 643) - (Step 644) + (Step 645) - (Step 646) + (Step 647) - (Step 648) + (Step 649) - (Step 650) + (Step 651) - (Step 652) + (Step 653) - (Step 654) + (Step 655) - (Step 656) + (Step 657) - (Step 658) + (Step 659) - (Step 660) + (Step 661) - (Step 662) + (Step 663) - (Step 664) + (Step 665) - (Step 666) + (Step 667) - (Step 668) + (Step 669) - (Step 670) + (Step 671) - (Step 672) + (Step 673) - (Step 674) + (Step 675) - (Step 676) + (Step 677) - (Step 678) + (Step 679) - (Step 680) + (Step 681) - (Step 682) + (Step 683) - (Step 684) + (Step 685) - (Step 686) + (Step 687) - (Step 688) + (Step 689) - (Step 690) + (Step 691) - (Step 692) + (Step 693) - (Step 694) + (Step 695) - (Step 696) + (Step 697) - (Step 698) + (Step 699) - (Step 700) + (Step 701) - (Step 702) + (Step 703) - (Step 704) + (Step 705) - (Step 706) + (Step 707) - (Step 708) + (Step 709) - (Step 710) + (Step 711) - (Step 712) + (Step 713) - (Step 714) + (Step 715) - (Step 716) + (Step 717) - (Step 718) + (Step 719) - (Step 720) + (Step 721) - (Step 722) + (Step 723) - (Step 724) + (Step 725) - (Step 726) + (Step 727) - (Step 728) + (Step 729) - (Step 730) + (Step 731) - (Step 732) + (Step 733) - (Step 734) + (Step 735) - (Step 736) + (Step 737) - (Step 738) + (Step 739) - (Step 740) + (Step 741) - (Step 742) + (Step 743) - (Step 744) + (Step 745) - (Step 746) + (Step 747) - (Step 748) + (Step 749) - (Step 750) + (Step 751) - (Step 752) + (Step 753) - (Step 754) + (Step 755) - (Step 756) + (Step 757) - (Step 758) + (Step 759) - (Step 760) + (Step 761) - (Step 762) + (Step 763) - (Step 764) + (Step 765) - (Step 766) + (Step 767) - (Step 768) + (Step 769) - (Step 770) + (Step 771) - (Step 772) + (Step 773) - (Step 774) + (Step 775) - (Step 776) + (Step 777) - (Step 778) + (Step 779) - (Step 780) + (Step 781) - (Step 782) + (Step 783) - (Step 784) + (Step 785) - (Step 786) + (Step 787) - (Step 788) + (Step 789) - (Step 790) + (Step 791) - (Step 792) + (Step 793) - (Step 794) + (Step 795) - (Step 796) + (Step 797) - (Step 798) + (Step 799) - (Step 800) + (Step 801) - (Step 802) + (Step 803) - (Step 804) + (Step 805) - (Step 806) + (Step 807) - (Step 808) + (Step 809) - (Step 810) + (Step 811) - (Step 812) + (Step 813) - (Step 814) + (Step 815) - (Step 816) + (Step 817) - (Step 818) + (Step 819) - (Step 820) + (Step 821) - (Step 822) + (Step 823) - (Step 824) + (Step 825) - (Step 826) + (Step 827) - (Step 828) + (Step 829) - (Step 830) + (Step 831) - (Step 832) + (Step 833) - (Step 834) + (Step 835) - (Step 836) + (Step 837) - (Step 838) + (Step 839) - (Step 840) + (Step 841) - (Step 842) + (Step 843) - (Step 844) + (Step 845) - (Step 846) + (Step 847) - (Step 848) + (Step 849) - (Step 850) + (Step 851) - (Step 852) + (Step 853) - (Step 854) + (Step 855) - (Step 856) + (Step 857) - (Step 858) + (Step 859) - (Step 860) + (Step 861) - (Step 862) + (Step 863) - (Step 864) + (Step 865) - (Step 866) + (Step 867) - (Step 868) + (Step 869) - (Step 870) + (Step 871) - (Step 872) + (Step 873) - (Step 874) + (Step 875) - (Step 876) + (Step 877) - (Step 878) + (Step 879) - (Step 880) + (Step 881) - (Step 882) + (Step 883) - (Step 884) + (Step 885) - (Step 886) + (Step 887) - (Step 888) + (Step 889) - (Step 890) + (Step 891) - (Step 892) + (Step 893) - (Step 894) + (Step 895) - (Step 896) + (Step 897) - (Step 898) + (Step 899) - (Step 900) + (Step 901) - (Step 902) + (Step 903) - (Step 904) + (Step 905) - (Step 906) + (Step 907) - (Step 908) + (Step 909) - (Step 910) + (Step 911) - (Step 912) + (Step 913) - (Step 914) + (Step 915) - (Step 916) + (Step 917) - (Step 918) + (Step 919) - (Step 920) + (Step 921) - (Step 922) + (Step 923) - (Step 924) + (Step 925) - (Step 926) + (Step 927) - (Step 928) + (Step 929) - (Step 930) + (Step 931) - (Step 932) + (Step 933) - (Step 934) + (Step 935) - (Step 936) + (Step 937) - (Step 938) + (Step 939) - (Step 940) + (Step 941) - (Step 942) + (Step 943) - (Step 944) + (Step 945) - (Step 946) + (Step 947) - (Step 948) + (Step 949) - (Step 950) + (Step 951) - (Step 952) + (Step 953) - (Step 954) + (Step 955) - (Step 956) + (Step 957) - (Step 958) + (Step 959) - (Step 960) + (Step 961) - (Step 962) + (Step 963) - (Step 964) + (Step 965) - (Step 966) + (Step 967) - (Step 968) + (Step 969) - (Step 970) + (Step 971) - (Step 972) + (Step 973) - (Step 974) + (Step 975) - (Step 976) + (Step 977) - (Step 978) + (Step 979) - (Step 980) + (Step 981) - (Step 982) + (Step 983) - (Step 984) + (Step 985) - (Step 986) + (Step 987) - (Step 988) + (Step 989) - (Step 990) + (Step 991) - (Step 992) + (Step 993) - (Step 994) + (Step 995) - (Step 996) + (Step 997) - (Step 998) + (Step 999) - (Step 1000) + (Step 1001) - (Step 1002) + (Step 1003) - (Step 1004) + (Step 1005) - (Step 1006) + (Step 1007) - (Step 1008) + (Step 1009) - (Step 1010) + (Step 1011) - (Step 1012) + (Step 1013) - (Step 1014) + (Step 1015) - (Step 1016) + (Step 1017) - (Step 1018) + (Step 1019) - (Step 1020) + (Step 1021) - (Step 1022) + (Step 1023) - (Step 1024) + (Step 1025) - (Step 1026) + (Step 1027) - (Step 1028) + (Step 1029) - (Step 1030) + (Step 1031) - (Step 1032) + (Step 1033) - (Step 1034) + (Step 1035) - (Step 1036) + (Step 1037) - (Step 1038) + (Step 1039) - (Step 1040) + (Step 1041) - (Step 1042) + (Step 1043) - (Step 1044) + (Step 1045) - (Step 1046) + (Step 1047) - (Step 1048) + (Step 1049) - (Step 1050) + (Step 1051) - (Step 1052) + (Step 1053) - (Step 1054) + (Step 1055) - (Step 1056) + (Step 1057) - (Step 1058) + (Step 1059) - (Step 1060) + (Step 1061) - (Step 1062) + (Step 1063) - (Step 1064) + (Step 1065) - (Step 1066) + (Step 1067) - (Step 1068) + (Step 1069) - (Step 1070) + (Step 1071) - (Step 1072) + (Step 1073) - (Step 1074) + (Step 1075) - (Step 1076) + (Step 1077) - (Step 1078) + (Step 1079) - (Step 1080) + (Step 1081) - (Step 1082) + (Step 1083) - (Step 1084) + (Step 1085) - (Step 1086) + (Step 1087) - (Step 1088) + (Step 1089) - (Step 1090) + (Step 1091) - (Step 1092) + (Step 1093) - (Step 1094) + (Step 1095) - (Step 1096) + (Step 1097) - (Step 1098) + (Step 1099) - (Step 1100) + (Step 1101) - (Step 1102) + (Step 1103) - (Step 1104) + (Step 1105) - (Step 1106) + (Step 1107) - (Step 1108) + (Step 1109) - (Step 1110) + (Step 1111) - (Step 1112) + (Step 1113) - (Step 1114) + (Step 1115) - (Step 1116) + (Step 1117) - (Step 1118) + (Step 1119) - (Step 1120) + (Step 1121) - (Step 1122) + (Step 1123) - (Step 1124) + (Step 1125) - (Step 1126) + (Step 1127) - (Step 1128) + (Step 1129) - (Step 1130) + (Step 1131) - (Step 1132) + (Step 1133) - (Step 1134) + (Step 1135) - (Step 1136) + (Step 1137) - (Step 1138) + (Step 1139) - (Step 1140) + (Step 1141) - (Step 1142) + (Step 1143) - (Step 1144) + (Step 1145) - (Step 1146) + (Step 1147) - (Step 1148) + (Step 1149) - (Step 1150) + (Step 1151) - (Step 1152) + (Step 1153) - (Step 1154) + (Step 1155) - (Step 1156) + (Step 1157) - (Step 1158) + (Step 1159) - (Step 1160) + (Step 1161) - (Step 1162) + (Step 1163) - (Step 1164) + (Step 1165) - (Step 1166) + (Step 1167) - (Step 1168) + (Step 1169) - (Step 1170) + (Step 1171) - (Step 1172) + (Step 1173) - (Step 1174) + (Step 1175) - (Step 1176) + (Step 1177) - (Step 1178) + (Step 1179) - (Step 1180) + (Step 1181) - (Step 1182) + (Step 1183) - (Step 1184) + (Step 1185) - (Step 1186) + (Step 1187) - (Step 1188) + (Step 1189) - (Step 1190) + (Step 1191) - (Step 1192) + (Step 1193) - (Step 1194) + (Step 1195) - (Step 1196) + (Step 1197) - (Step 1198) + (Step 1199) - (Step 1200) + (Step 1201) - (Step 1202) + (Step 1203) - (Step 1204) + (Step 1205) - (Step 1206) + (Step 1207) - (Step 1208) + (Step 1209) - (Step 1210) + (Step 1211) - (Step 1212) + (Step 1213) - (Step 1214) + (Step 1215) - (Step 1216) + (Step 1217) - (Step 1218) + (Step 1219) - (Step 1220) + (Step 1221) - (Step 1222) + (Step 1223) - (Step 1224) + (Step 1225) - (Step 1226) + (Step 1227) - (Step 1228) + (Step 1229) - (Step 1230) + (Step 1231) - (Step 1232) + (Step 1233) - (Step 1234) + (Step 1235) - (Step 1236) + (Step 1237) - (Step 1238) + (Step 1239) - (Step 1240) + (Step 1241) - (Step 1242) + (Step 1243) - (Step 1244) + (Step 1245) - (Step 1246) + (Step 1247) - (Step 1248) + (Step 1249) - (Step 1250) + (Step 1251) - (Step 1252) + (Step 1253) - (Step 1254) + (Step 1255) - (Step 1256) + (Step 1257) - (Step 1258) + (Step 1259) - (Step 1260) + (Step 1261) - (Step 1262) + (Step 1263) - (Step 1264) + (Step 1265) - (Step 1266) + (Step 1267) - (Step 1268) + (Step 1269) - (Step 1270) + (Step 1271) - (Step 1272) + (Step 1273) - (Step 1274) + (Step 1275) - (Step 1276) + (Step 1277) - (Step 1278) + (Step 1279) - (Step 1280) + (Step 1281</p>

13. Zhang, W.; Xu, H.; Wang, Z.; He, Z.; Zhu, Z.; Ren, K. Can Small Language Models Reliably Resist Jailbreak Attacks? A Comprehensive Evaluation. *arXiv 2025*, 2503.06519. DOI: 10.48550/arXiv.2503.06519.
14. Shim, J.; Dong Ho, S.; Kim, J. CoT Harms Performance of Rather Smaller Language Models. *Proceedings of the International Conference on Industrial Engineering and Operations Management (1st World Congress) 2024*, 504–508. DOI: 10.46254/WC01.20240169.
15. Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; Sutskever, I. Language Models Are Unsupervised Multitask Learners. *OpenAI blog 2019*
16. Mathematical Association of America. American Mathematics Competitions. MAA Website. <https://maa.org/student-programs/amc/> (accessed 2025-07-24).
17. Zhang, J.; Zhao, Y.; Zhang, L.; Hu, J.; Luan, X.; Xu, Z.; Yang, F. Psychometric-Based Evaluation for Theorem Proving with Large Language Models. *arXiv 2025*, abs/2502.00855. DOI: 10.48550/arXiv.2502.00855.
18. Google DeepMind. Gemma 3 model card. <https://huggingface.co/google/gemma-3-27b-it> (accessed 2025-06-30).
19. Atil, B.; Aykent, S.; Chittams, A.; Fu, L.; Passonneau, R. J.; Radcliffe, E.; Rajagopal, G. R.; Sloan, A.; Tudrej, T.; Ture, F.; Wu, Z.; Xu, L.; Baldwin, B. Non-Determinism of “Deterministic” LLM Settings. *arXiv 2024*, 2408.04667. DOI: 10.48550/arXiv.2408.04667.
20. Ahn, J.; Yin, W. Prompt-Reverse Inconsistency: LLM Self-Inconsistency Beyond Generative Randomness and Prompt Paraphrasing. *arXiv 2025*, 2504.01282. DOI: 10.48550/arXiv.2504.01282.
21. Kojima, T.; Gu, S. S.; Reid, M.; Matsuo, Y.; Iwasawa, Y. Large Language Models Are Zero-Shot Reasoners. *arXiv 2022*, 2205.11916. DOI: 10.48550/arXiv.2205.11916.
22. Ji, Z.; Yu, T.; Xu, Y.; Lee, N.; Ishii, E.; Fung, P. Towards Mitigating Hallucination in Large Language Models via Self-Reflection. Findings of the Association for Computational Linguistics: *EMNLP 2023 2023*, Findings, 1827–1843. DOI: 10.18653/v1/2023.findings-emnlp.123.
23. Shi, F.; Chen, X.; Misra, K.; Scales, N.; Dohan, D.; Chi, E. H.; Schärli, N.; Zhou, D. Large Language Models Can Be Easily Distracted by Irrelevant Context. *Proceedings of the 40th International Conference on Machine Learning 2023*, 202, 31210–31227
24. Levy, M.; Jacoby, A.; Goldberg, Y. Same Task, More Tokens: the Impact of Input Length on the Reasoning Performance of Large Language Models. Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL 2024), Long Papers 2024, 1 (Long Papers), 15339–15353. DOI: 10.18653/v1/2024.acl-long.818.
25. Zhang, Z.; Li, J.; Lan, Y.; Wang, X.; Wang, H. An Empirical Study on Prompt Compression for Large Language Models. *Building Trust Workshop at ICLR 2025*, 2025. DOI: 10.48550/arXiv.2505.00019.
26. Schreiter, D. Prompt Engineering: How Prompt Vocabulary Affects Domain Knowledge. *arXiv 2025*, 2505.17037. DOI: 10.48550/arXiv.2505.17037.
27. Cao, B.; Cai, D.; Zhang, Z.; Zou, Y.; Lam, W. On the Worst Prompt Performance of Large Language Models. *arXiv 2024*, 2406.10248. DOI: 10.48550/arXiv.2406.10248.
28. Han, D.; Xia, M.; Madrigal Diaz, D.; Kessler, S.; Mallick, A.; Zhang, X.; Garcia, M. D. C. H.; Xu, J.; Rühle, V.; Rajmohan, S. Enhancing Reasoning Capabilities of Small Language Models with Blueprints and Prompt Template Search. *arXiv 2025*, 2506.08669. DOI: 10.48550/arXiv.2506.08669.
29. Kim, Y.; Yi, E.; Kim, M.; Yun, S. Y.; Kim, T. Guiding Reasoning in Small Language Models with LLM Assistance. *arXiv 2025*, 2504.09923. DOI: 10.48550/arXiv.2504.09923.
30. Qi, Z.; Ma, M.; Xu, J.; Zhang, L. L.; Yang, F.; Yang, M. Mutual Reasoning Makes Smaller LLMs Stronger Problem-Solvers. *arXiv 2024*, 2408.06195. DOI: 10.48550/arXiv.2408.06195.
31. Liu, X.; Xu, P.; Wu, J.; Yuan, J.; Yang, Y.; Zhou, Y.; Liu, F.; Guan, T.; Wang, H.; Yu, T.; McAuley, J.; Ai, W.; Huang, F. Large Language Models and Causal Inference in Collaboration: A Comprehensive Survey. *arXiv 2024*, 2403.09606. DOI: 10.48550/arXiv.2403.09606.

## ■ Author

Eric is currently a high school student in the United States. He is very passionate about both artificial intelligence and competitive mathematics, and has a strong background in competitive math, having qualified for the USAMO in 2025, and USAJMO in both 2023 and 2024.