

The Illusion of Diminishing Returns: Measuring Long Horizon Execution in LLMs

Akshit Sinha^{1*} Arvindh Arun^{2*} Shashwat Goel^{3,4*}
Steffen Staab^{2,5} Jonas Geiping^{3,4,6}

¹University of Cambridge ²Institute for AI, University of Stuttgart

³Max Planck Institute for Intelligent Systems ⁴ELLIS Institute Tübingen

⁵University of Southampton ⁶Tübingen AI Center



Code



Dataset

Abstract

Does continued scaling of large language models (LLMs) yield diminishing returns? Real-world value often stems from the length of task an agent can complete. We start this work by observing the simple but counterintuitive fact that marginal gains in single-step accuracy can compound into exponential improvements in the length of a task a model can successfully complete. Then, we argue that failures of LLMs when simple tasks are made longer arise from mistakes in *execution*, rather than an inability to *reason*. We propose isolating *execution* capability, by explicitly providing the *knowledge* and *plan* needed to solve a long-horizon task. We find that larger models can correctly execute significantly more turns even when small models have 100% single-turn accuracy. We observe that the per-step accuracy of models degrades as the number of steps increases. This is not just due to long-context limitations—curiously, we observe a *self-conditioning* effect—models become more likely to make mistakes when the context contains their errors from prior turns. Self-conditioning does not reduce by just scaling the model size. In contrast, recent *thinking* models do not self-condition, and can also execute much longer tasks in a single turn. We conclude by benchmarking frontier thinking models on the length of task they can execute in a single turn. Overall, by focusing on the ability to execute, we hope to reconcile debates on how LLMs can solve complex reasoning problems yet fail at simple tasks when made longer, and highlight the massive benefits of scaling model size and sequential test-time compute for long-horizon tasks.

1 Introduction

Is continued scaling of compute for Large Language Models (LLMs) economically justified given diminishing marginal gains? This question lies at the heart of the ongoing debate on the viability of continued massive investments in LLMs. While scaling laws show diminishing returns on metrics like test loss, the economic potential of LLMs might arise from automating long, multi-step tasks (METR, 2025). However, long-horizon tasks have been the Achilles’ heel of Deep Learning. We saw impressive self-driving demos take over a decade to translate to reliability in long-distance driving. Vision models can generate impressive images, and yet consistency over long videos remains an unsolved challenge. As the industry races to build agents that tackle entire projects, not just isolated questions, a fundamental question arises: *How can we measure the number of steps an LLM can reliably execute?*

LLM failures on simple, but long tasks have been considered a fundamental inability to *reason* (Mirzadeh et al., 2024). Despite massive improvements on complex reasoning benchmarks, Shojaei et al. (2025) claim *thinking models* (Guo et al., 2025) only give an “illusion of thinking”, as they eventually fail when the task is made longer. These results have sparked much debate in the community, which we think can be resolved by decoupling the need for *planning*

*Equal contribution

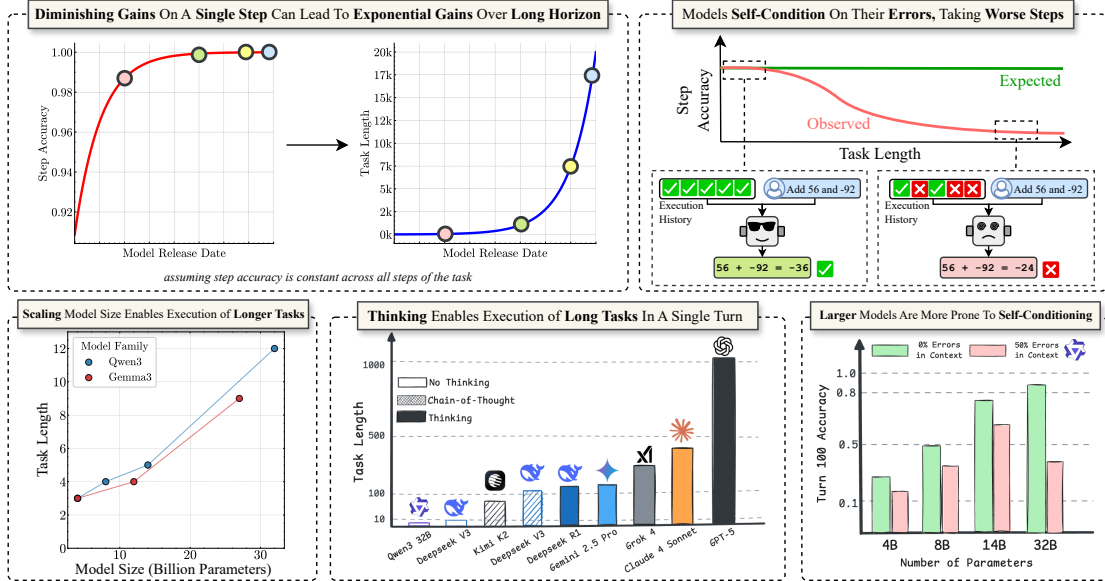


Figure 1: A Summary of our contributions. We note that diminishing returns can enable exponential gains in horizon length (that is the length of tasks a model can complete). We design a simple task that isolates the capability of long-horizon execution in language models, ablating the need for knowledge and planning. We find that frontier models benefit considerably from both scaling model size and test-time compute when executing long horizon tasks.

and *execution* in reasoning or agentic tasks. *Planning* involves deciding what information to retrieve or tools to use and in which order, while *execution* involves carrying out the plan. In [Shojaei et al. \(2025\)](#), the LLM clearly knows the plan, as it initially follows it for many steps correctly. We posit that the eventual failures are in execution—as the task gets longer, the model is more likely to make a mistake in executing the plan. Although much attention has been paid to LLM planning abilities ([Kambhampati et al., 2024](#)), execution remains an understudied challenge, despite being increasingly important as LLMs begin to be used for long reasoning and agentic tasks.

In this work, we measure *long-horizon execution* capabilities of LLMs in a controlled setting. We isolate the *execution* capability of LLMs by explicitly providing them the *knowledge* and *plan* needed. By controlling the number of turns, and the number of steps per turn, which together contribute to task length, we reveal insights about long-horizon execution in LLMs:

Does Scaling have Diminishing Returns? We observe that diminishing improvements in single-step accuracy can compound, leading to exponential growth in the length of task a model can complete. Traditionally, scaling model size is assumed to increase capacity to store parametric knowledge or search for plans. Yet, even when the required knowledge and plan are explicitly provided, empirically we find that scaling model size leads to large improvements in the number of turns a model can execute successfully.

The Self-Conditioning Effect. One might assume that failures on long tasks are simply due to the compounding of a small, constant per-step error rate. However, we find that the per-step error rate itself rises as the task progresses. This is in contrast to humans, who typically improve at executing a task with practice. We hypothesize that as a significant fraction of model training is to predict the most likely next token given its context, conditioning models on their own error-prone history increases the likelihood of future errors. We test this by controlling the error rate in the history shown to the model. As the error rate in the history is increased, we observe a sharp degradation in subsequent step accuracy, validating that models *self-condition*. We show self-conditioning leads to degradation in model performance in long-horizon tasks beyond previously identified long-context issues, and unlike the latter, is not mitigated by scaling model size.

The Impact of Thinking. We find recent thinking models are not affected by prior mistakes, fixing self-conditioning. Further, sequential test time compute greatly improves the length of task a model can complete in a single turn. Where without CoT, frontier LLMs like DeepSeek-

V3 fail at performing even two steps of execution, its thinking version R1 can execute 200, highlighting the importance of reasoning before acting (Yao et al., 2023). We benchmark frontier thinking models, and find GPT-5 thinking (codename “Horizon”) can execute over 1000 steps, far ahead of the next best competitor, Claude-4-Sonnet at 432.

The “jagged frontier” (Dell’Acqua et al., 2023) of LLM capabilities remains fascinating yet confusing. Unlike traditional machines, LLMs are more susceptible to failure when used for executing repetitive tasks. Thus, we argue execution failures in long tasks should not be misinterpreted as the inability to reason or plan. We show long-horizon execution improves dramatically by scaling model size and sequential test time compute. If the length of tasks a model can complete indicates its economic value, continued investment in scaling compute might be worth the cost, even if short-task benchmarks give the illusion of slowing progress.

2 Formulation

In an agentic or reasoning task, the model begins in an initial state (based on the first input) and has to perform a sequence of steps to reach the final goal. A long-horizon task requires a large number of steps, where the task length is the number of steps needed to complete it. We define the following metrics to evaluate performance:

Step Accuracy. It measures the fraction of samples where the state update from step $i - 1$ to step i is correct, regardless of the correctness of the model’s state at step $i - 1$.

Turn Accuracy. A turn is a single interaction with the model, which may require executing multiple steps. Turn Accuracy measures the fraction of samples where the state update from turn $t - 1$ to turn t is correct, regardless of the correctness of the model’s state at turn $t - 1$.

Turn Complexity (K). It is defined as the number of steps the model has to execute per turn.

Task Accuracy. It measures the fraction of samples in which the model can complete a task of i steps without making any mistakes in the process.

Horizon Length (H_s). We define the horizon length of a model given a success rate threshold $0 \leq s \leq 1$ as the first step i where the model’s mean task accuracy across samples drops below s . It can be interpreted as: the model can perform a task of length H_s without making mistakes, with probability s . We use $s = 0.5$ unless otherwise specified, in analogy to Kwa et al. (2025).

2.1 Diminishing returns in Step Accuracy yield exponential gains on Horizon Length

We begin by analyzing the relationship between a model’s single-step accuracy and its horizon length. To obtain a mathematical relation, we make two simplifying assumptions similar to LeCun (2023). First, we assume a model’s step accuracy remains constant over the task. Second, we assume a model does not self-correct, meaning any single error leads to task failure. We assume this only for the analysis here, which is illustrative and provides useful intuition. Our empirical analysis goes beyond this, investigating how LLMs, in fact, do not exhibit constant step accuracy for long horizon execution, and may correct mistakes.

Proposition 1. Assuming a constant step accuracy p and no self-correction, the task-length H at which a model achieves a success rate s is given by:

$$H_s(p) = \left\lceil \frac{\ln(s)}{\ln(p)} \right\rceil \approx \frac{\ln(s)}{\ln(p)}$$

(The derivation is provided in Section H.)

We plot this growth function in Figure 2 for $s = 0.5$. Note how after the step accuracy crosses 70%, small gains in step accuracy lead to faster than exponential improvement in horizon

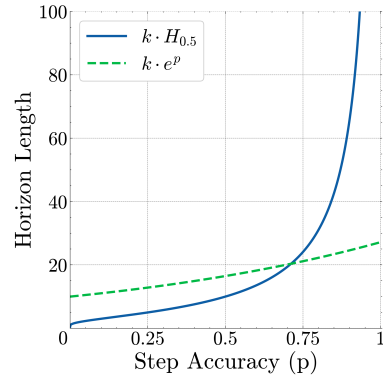


Figure 2: Growth of Horizon Length. The length of task a model can perform at more than 50% accuracy grows faster than exponential as a function of the step accuracy after the 70% mark.

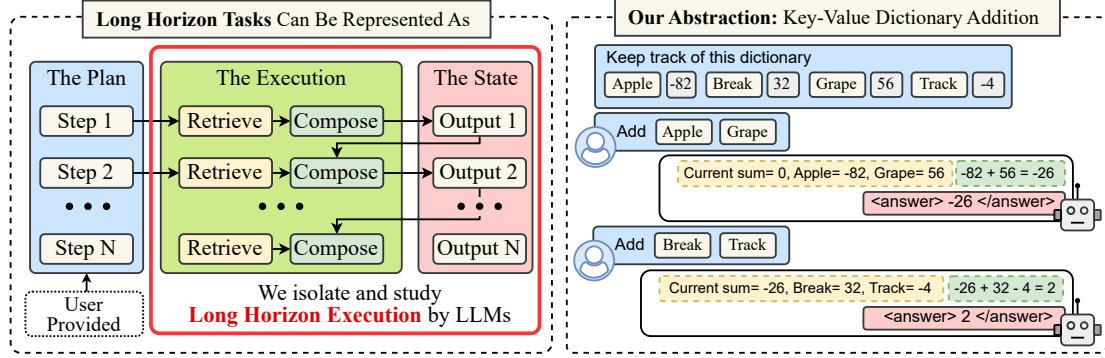


Figure 3: Overview of our framework. (Left) Our framework models long-horizon tasks as a sequence of *retrieve-then-compose* steps. (Right) We design a simple task that decouples planning from execution: in each turn, we provide the model the plan as key(s), asking it to *retrieve* their value(s), and *compose* them to maintain a running sum.

length. This derivation shows that even if accuracy improvements seem to slow down on question answering benchmarks, which typically have short tasks (METR, 2025), one can still mathematically expect large gains on longer tasks.

For example, on software engineering tasks, Kwa et al. (2025) empirically observed that the horizon length at $s = 0.5$ of frontier models is growing exponentially, doubling every 7 months. Using our result above, in Figure 1 we show that such exponential growth in horizon length occurs even in a regime of diminishing returns on step accuracy. If we set $s = 0.5$, we obtain $H_{0.5} = -\frac{\ln(2)}{\ln(p)}$. As such, the step-accuracy p required to sustain exponential growth in $H_{0.5}$ over time (x) is $2^{\frac{-1}{2^x}}$, which is indeed a diminishing function.

We note that human labor is often compensated for its time. If the economic value of an agent also arises from the length of tasks it can complete, single-turn or short task benchmarks may be an illusory reference for evaluating the benefits of further investment in LLM compute. They might give a mirage of slowing progress, while the length of tasks a model can complete, which we think is a better indicator of economic value, continues to grow fast.

2.2 Isolating execution by decoupling planning and knowledge

We now describe how we measure long-horizon execution empirically. First, as a motivating example, consider an agent for the popular, economically valuable, task of booking flights. Upon receiving a search result, it must evaluate the surfaced flights to determine which one to book. The *plan* for assessing a single flight option may involve a sequence of actions, such as viewing detailed information, verifying that the flight timings, baggage allowance, and airline reviews align with user preferences, applying any available discounts or reward programs, and ultimately making a selection based on cost and travel time. Each of these individual steps requires *retrieving* some information, and *composing* it with the existing information state to eventually evaluate one flight option, and both of these operations require *knowledge*. The successful evaluation of multiple flight options constitutes the *execution* of this plan until a final booking decision is made.

In this work, we focus on *execution*, as we argue that it is a critical component of long-horizon capabilities. Execution has traditionally received less attention than capabilities such as reasoning, planning, and world knowledge, which have been the primary focus of LLM capability discussions. This relative neglect is significant to the extent that failures in execution have been misattributed to limitations in reasoning or planning capabilities (Shojaee et al., 2025; Khan et al., 2025). This perception may stem from the view that execution is a straightforward or mundane task. After all, this is what machines have been historically good at. Humans, too, are quite reliable at executing a task once they learn how to do it, even improving with practice. However, as LLMs do not come with correctness guarantees, we posit that execution can be surprisingly challenging for an LLM over a long horizon. We hypothesize that:

*Even if reasoning, planning, and world knowledge are perfected,
LLMs will still make mistakes in execution over a long-horizon.*

To demonstrate this, we isolate execution failures by explicitly providing the requisite knowledge and plan. We chain the *retrieve-then-compose* step motivated in the flight-selection agent example above. Each step involves *retrieving* relevant information or a tool specified by the plan and then *composing* its output to update the current state. The plan is deciding what to retrieve and how to compose it, whereas execution is actually performing those operations. This fits a natural abstraction—a key-value dictionary. The *key* serves as one step of a plan specifying what knowledge to retrieve, or tool to call, while the *value* represents the knowledge or tool output, which then has to be composed with the current state. In our study, we provide the plan as the keys in each query, eliminating the need for *planning* abilities from the LLM. We also provide the key-value dictionary in context, removing any dependency on the model’s parametric *knowledge*. With this design, we directly control two important axes that multiply to obtain the task length (number of retrieve-then-compose steps): the number of turns, and the turn complexity (K). The turn complexity can be varied by changing the number of keys queried per turn.

3 Experiments

Setup. As illustrated in Figure 3, we provide the model with the needed *knowledge*, a fixed, in-context dictionary $\mathcal{D} : \mathcal{V} \rightarrow \mathbb{Z}$, where \mathcal{V} is a vocabulary of common five-letter English words and values are integers sampled uniformly from $[-99, 99]$. The initial state is $S_0 = 0$. In turn $t \in \{1, \dots, T\}$, the model receives an explicit *plan* $P_t = \{k_{t,1}, \dots, k_{t,K}\}$, which is a set of K keys sampled from \mathcal{V} . For each turn t , the model must execute this plan, which requires updating the state, S_t to maintain a running sum of values for all past queried keys. This requires the retrieve-then-compose steps defined above:

1. **Retrieval:** Look up the integer value $\mathcal{D}[k]$ for each key $k \in P_t$
2. **Composition:** Sum these values and add them to the previous state, $S_t = S_{t-1} + \sum_{i=1}^K \mathcal{D}[k_{t,i}]$

We choose short English words and two-digit integers to minimize errors arising from tokenization. The state transition here is Markovian, depending only on S_{t-1} and P_t . The task is extremely simple, by design, to isolate long-horizon execution by minimizing the knowledge needed for the retrieval and composition operation. More details, including the exact prompt, are provided in Section E. We analyze empirical performance on the individual retrieval, and composition operations in Section D and the format following errors in Section F.

3.1 Effect of increasing the number of turns

We first test our hypothesis that long-horizon execution can be challenging even on tasks where world-knowledge and planning are not required. We then study the benefits of scaling model size on long-horizon execution.

Setup. We evaluate the Qwen3 (Yang et al., 2025) and Gemma3 (Gemma-Team et al., 2025) model families, as they offer a range of sizes: [4, 8, 14, 32]B and [4, 12, 27]B parameters, respectively. For this experiment, we set the turn complexity to its simplest form ($K = 1$), providing a single key per turn, and vary the number of turns. Models are instructed to output the final answer directly, without intermediate thinking tokens, with the format enforced via few-shot examples. In Section C we show the results below also hold with chain of thought and thinking models.

Result 1: Execution Alone is Challenging. We present the results in Figure 4. All models except Gemma3-4B and Qwen3-4B achieve 100% accuracy on the first step, highlighting how they have the knowledge and reasoning capability required to perfectly do a single step of our task. Yet, task accuracy falls rapidly over subsequent turns. Even the best-performing model (Qwen3-32B) sees its accuracy fall below 50% within 15 turns. This confirms our hypothesis that long-horizon execution can be challenging for LLMs even when planning and knowledge requirements are removed.

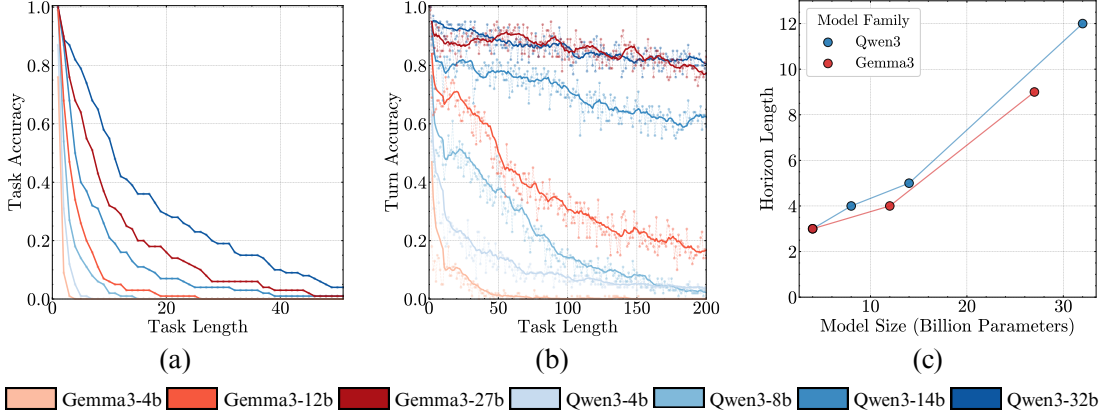


Figure 4: Scaling model size non-diminishingly improves the number of turns it can execute. We vary the model size and study both full task (a) and turn-wise accuracy (b) as the number of turns increases. Bold lines are a running average of accuracy over 5 turns. The dotted lines (turn-wise accuracy) in (b) show single-step accuracy for our task is 100% for all except the smallest models. Yet, as the number of turns increases, the performance gap between small and large models widens (a), with the latter having significantly more horizon length (c).

Result 2: Non-Diminishing Benefit of Scaling Model Size. As shown in Figure 4 (a), larger models sustain higher task accuracy for significantly more turns, resulting in a clear scaling trend for horizon length (Figure 4 (c)). We abstain from deriving a “scaling law” since we can only obtain at most four model sizes from the same family, but the improvements do not seem diminishing. This observation is non-trivial. While the benefits of increasing model size are often attributed to improved knowledge capacity, our task is not knowledge-constrained, as even small models achieve perfect single-step accuracy, nor is the task more complex so that a larger model would be required. Yet, larger models are clearly more reliable at executing the task for longer. A possible explanation is the redundancy of internal circuits in larger models, which ensembles to reduce error (Lindsey et al., 2025). However, we find that simulating this redundancy with output-level aggregation of parallel compute (Section B) does not replicate the gains observed from scaling model size.

Takeaway 1. Long-horizon execution is challenging. Scaling model size significantly increases the number of turns a model can correctly execute.

3.2 Why Does Turn Accuracy Degrade? The Self-Conditioning Effect

One might expect a model’s per turn performance to remain constant. Yet, Figure 4(b) shows the accuracy of individual turns steadily degrades as the number of turns increases. We investigate two competing hypotheses:

- 1. Degradation as the context length increases.** The model’s performance degrades simply due to increasing context length (Zhou et al., 2025a), irrespective of its content.
- 2. Self-conditioning.** The model conditions on its own past mistakes. It becomes more likely to make a mistake after observing its own past errors in previous turns.

Setup. To disentangle these factors, we conduct a counterfactual experiment by manipulating the model’s chat history. We control the error rate by injecting artificial output histories with a chosen error rate in the same format. If we fully *heal* the history, with a 0% error rate, degradation in the model’s turn accuracy between turn 1 and a later turn can be attributed to long-context issues. If a model’s accuracy for a fixed later turn consistently worsens with increasing error rate in prior turns, this would demonstrate that models condition on their past mistakes, increasing the likelihood of future errors.

Result 3: Self-Conditioning causes degradation in turn accuracy beyond long-context. Our results in Figure 5 (a) show evidence for degradation due to both long-context and self-

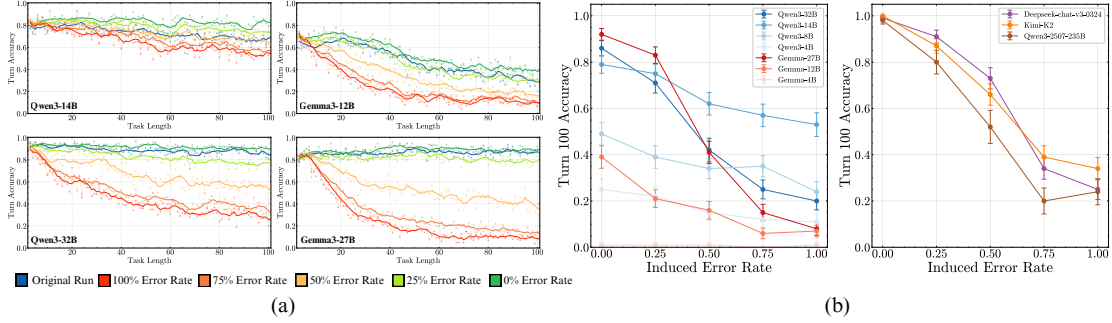


Figure 5: Language models self-condition on their previous mistakes, leading to more mistakes in subsequent turns. By manipulating the chat history, we counterfactually vary the fraction of errors in previous turns. We find this increases the likelihood of errors in future turns (left). This shows a source of degradation in turn-wise model accuracy beyond long-context, as in the turn 100 slice (right) model accuracies are much higher when we provide a fully correct history. Scaling model size increases self-conditioning, even for frontier non-thinking models.

conditioning. When conditioned on an error-free history (Induced Error Rate = 0.00), model turn accuracy at turn 100 is below its initial value, consistent with prior observations of long-context degradation (Zhou et al., 2025a). More interestingly, as we increase the rate of injected errors into the context, accuracy at turn 100 consistently degrades further. This demonstrates the self-conditioning effect—as models make mistakes, they become more likely to make more mistakes, leading to a continuous degradation in per-turn accuracy throughout the output trajectory as shown in Figure 5 (b).

Result 4: Unlike long-context, scaling model size does not mitigate self-conditioning. Notice that the accuracy at turn 100 at the induced error rate of 0 consistently improves for larger models. As shown in Figure 5 (c), scaling to frontier (200B+ parameter) models like Kimi-K2 (Kimi-Team et al., 2025), DeepSeek-V3 (DeepSeek-AI et al., 2025), and Qwen3-235B-Instruct-2507 (Yang et al., 2025) largely solves long-context degradation for up to 100 turns, achieving near-perfect accuracy on a healed history. However, even these large models remain susceptible to self-conditioning, as their performance consistently degrades as the induced error rate in their history increases. This may be akin to recent results showing larger models shift in personality during multi-turn conversations (Choi et al., 2024; Becker et al., 2025), where in our case, the drift is toward a personality that makes errors.

Takeaway 2. Models *self-condition* on their previous mistakes, leading to degradation in per-step accuracy. Scaling model size is not sufficient to mitigate this.

We now study the effect of enabling sequential test time compute (“thinking”) for these models.

Setup. We enable thinking for the Qwen3 models, which are post-trained with reinforcement learning (RL). These models are trained to generate reasoning traces even when the context contains only the final answers from previous turns. This contrasts with standard chain-of-thought prompting, where models often fail to reason if prior reasoning steps are omitted from the context. We found that the Gemma3 models, when prompted for CoT, were unable to generate reasoning if prior traces were omitted from the context. These models exhibited a form of format-based self-conditioning: after observing a history of turns with only final outputs, they would ignore explicit user instructions to think step by step and

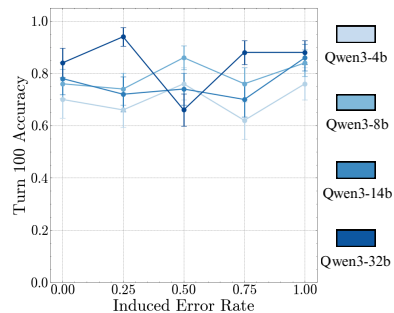


Figure 6: Thinking fixes self-conditioning. Qwen3 models with thinking enabled no longer self-condition, even when the entire prior history has wrong answers, in contrast to non-thinking results.

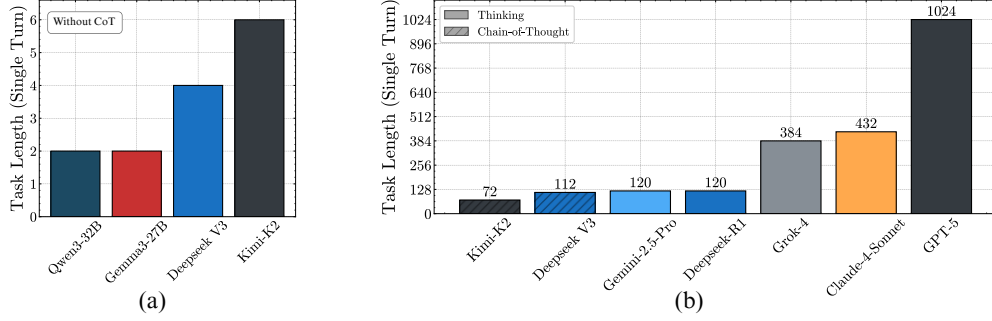


Figure 7: Benchmarking the length of task models can execute in a single turn. Without CoT or thinking, even the biggest models fail to execute more than a few steps (left). Sequential test time compute (thinking tokens) significantly improve this, especially when trained with RL (eg R1 vs DeepSeek V3) (right). GPT-5 is far ahead of the rest, executing over 1000 steps, with Claude-4-Sonnet second at around 400.

revert to producing only a final answer, further discussed in Section G. Given this limitation of CoT prompted Gemma3 models, we just focus on studying the self-condition effect on thinking-enabled Qwen3 models, by observing turn 100 accuracy while controlling the error rate in prior turns as before.

Result 5: Thinking Fixes Self-Conditioning. In Figure 6, we clearly find that the Qwen3 thinking models do not self-condition—the models’ accuracy at turn 100 remains stable, regardless of the error rate in its context. This could arise from two reasons. First, RL training can reduce the most likely next token prediction behaviour of language models, making them oriented towards task success rather than continuing the context. Second, the removal of thinking traces from prior turns could reduce the influence of prior turns on the model’s output, as it thinks about the new turn independently. By inspecting the models’ thinking traces, we observe that they do not refer back to prior turns in their chain of thought. Furthermore, we experiment with context management by explicitly removing prior history as a potential fix, and find that it indeed mitigates self-conditioning (Section A).

3.3 What is the length of tasks models can complete in a single turn?

In the previous sections, we measured how many turns models can successfully execute a single retrieve-then-compose step. However, most real-world tasks require more complex processing every turn. The total task length a model can handle is a function of both the number of turns and the number of steps to execute per turn. We now measure the latter dimension: the maximum number of steps a model can execute per turn.

Setup. To quantify this capability, we propose a benchmark that measures the maximum single-turn execution capacity of various models. We run a binary search (Lehmer, 1960) to find the highest turn complexity (K , the number of keys) the model can provide the correct sum for with accuracy $\geq 80\%$. We evaluate a suite of frontier models like GPT-5 (OpenAI, 2025), Claude-4 Sonnet (Anthropic, 2025), Grok 4 (xAI, 2025), Gemini 2.5 Pro (Gemini Team, 2025), Kimi K2 (Kimi-Team et al., 2025), Qwen3-Instruct-235B-2507 (Yang et al., 2025), and DeepSeek R1 (Guo et al., 2025). An advantage of our benchmark is that it is contamination-free, as new examples can be generated programmatically.

Result 1: Without chain of thought, non-thinking models struggle to chain even two steps in a single turn. In Figure 12 (left), we first find that when prompted to directly answer, without chain-of-thought, the larger Qwen3 32B, Gemma3 27B, as well as frontier non-think models like DeepSeek-V3 (670B), and Kimi K2 (1026B), fail to execute even a turn complexity of 2. This is consistent with prior work showing the necessity of thinking tokens for transformers to perform sequential tasks (Weiss et al., 2021; Merrill and Sabharwal, 2023). We highlight this because many agentic workflows directly ask the model to act, without chain-of-thought, to fit more actions in the context window. We see that the number of steps the model can execute in a single turn improves significantly with chain-of-thought. It shows the importance of reasoning before

acting (ReAct (Yao et al., 2023)) for agents. In Section B, we also show that parallel test time compute like majority voting (Wang et al., 2022) only leads to marginal improvements in both single turn execution length and number of turns. This provides preliminary evidence that for long-horizon execution, sequential test-time compute is more effective.

Result 2: Benchmarking Frontier Models. In Figure 12 (right), we benchmark frontier models on the length of task they can execute in a single turn. We find a surprisingly large gap between GPT-5 (codenamed Horizon) and others like Gemini 2.5 Pro, Grok 4, and DeepSeek R1. We also find that the RL-trained thinking model DeepSeek R1 substantially outperforms its instruction-tuned counterpart, DeepSeek-V3. Overall, long-horizon execution is a challenge in which open-weight models are still catching up to those only available through APIs, highlighting opportunities for future research.

Takeaway 3. Thinking models fix self-conditioning, and can also execute significantly longer tasks in a single turn.

4 Related Work

Long Context. Much prior work has focused on improving the maximum context length that can be provided in the input to a language model (Su et al., 2021), and evaluating whether (Tay et al., 2020) and how (Olsson et al., 2022; Li et al., 2023) models maintain performance as the context gets longer (Tay et al., 2020). Closest is the recent RULER (Hsieh et al., 2024) and GSM-Infinite (Zhou et al., 2025b), which also uses synthetic data to systematically evaluate long-context abilities. While long-context will help models execute for longer, it is a different capability compared to long-horizon execution (Zhou et al., 2023; Chen et al., 2024a), as it focuses on performance as a function of input, not output length. We identified one such difference, the self-conditioning effect—where past errors in model output increase the chance of future mistakes, and disentangle this effect from long-context degradation in Section 3.2.

Classical Reasoning and Planning. Automated planning, especially when formulated for discrete and deterministic spaces, has long been a mainstay of artificial intelligence research. In symbolic AI, once tasks are formalized, for example into STRIPS plans (Fikes and Nilsson, 1971), they can be evaluated without issues in execution. Prior work (Chen et al., 2024b; Valmeekam et al., 2024) has shown LLMs struggle to match symbolic algorithms for automated planning. In contrast, we focus on straightforward execution of provided plans over a long horizon.

Increasing Task Complexity (length). Multiple works have recently shown how models worsen as “problem complexity” increases Zhou et al. (2025b), often attributed to failures of reasoning (Cheng, 2025; Shojaei et al., 2025). Recently, multiple real-world long-horizon agentic benchmarks have been proposed (Backlund and Petersson, 2025; Xie et al., 2024; Shen et al., 2025), where prior work has studied planning failures (Chen et al., 2024b). By designing a task where no reasoning is required, given that we provide the model the requisite plan and knowledge, we show that execution alone can be a challenge (Zhu et al., 2025; Sun et al., 2025), degrading model accuracy on longer tasks. Our observations on scaling could hold for the related problem of length-generalization—training models to succeed on tasks longer than those seen during training (Fan et al., 2024; Cai et al., 2025).

Controlled Evaluations with Synthetic Data. Improvements on real-world benchmarks are the ultimate measure of AI progress, but understanding LLM capabilities and shortcomings sometimes requires disentangling the many factors that compound in real tasks. Our empirical approach of performing a controlled study of LLM capabilities, using a simplified task to remove confounders, aligns with recent work on architecture design (Allen-Zhu, 2024; Poli et al., 2024), recall from parametric memory (Arora et al., 2023) (where our retrieve step requires in-context retrieval), length generalization (Lee et al., 2025), and the ability to form new abstractions (Chollet et al., 2024). We focus on a different capability—long-horizon execution—which we posit is becoming increasingly important as we enter the era of experience (Silver and Sutton, 2025).

5 Discussion

Scaling laws for language models show diminishing returns on the loss for the single step of predicting the next token (Kaplan et al., 2020; Hoffmann et al., 2022). When models competed in simple knowledge-based question-answering tasks such as MMLU (Hendrycks et al., 2020), such single-step measurements could inform us about the rate of progress. This has changed in the last year. Where earlier we could only post-train on human demonstrations (Mishra et al., 2021), language models can now be trained with just rewards (Shao et al., 2024), enabling sophisticated reasoning (Guo et al., 2025; Jain et al., 2024) and agents (Kimi Team et al., 2025). This opens up the opportunity to solve much longer tasks where earlier human supervision would be too expensive to scale. Our work shows how diminishing returns on single-step performance can compound to provide large benefits in the length of tasks a model can solve. This motivates the need to study empirical scaling laws for horizon length in agents (Hilton et al., 2023). An astute reader might wonder if execution failures should be solved by providing the model access to tools (Schick et al., 2023). Tools indeed help shift the burden of execution from probabilistic models to reliable programs. However, reasoning is often fuzzy, and not always easy to implement as a tool, requiring the model to execute some by itself. Even calling the right tools requires reliable execution from the model (Patil et al., 2025).

While there has been recent interest in evaluations for LLM reliability (Vendrow et al., 2025; Yao et al., 2024), they do not focus on long-horizon outputs, where the context differs every step. For example, Yao et al. (2024) focus on the $pass^k$ metric, which checks if the model makes a mistake when we sample k generations. However, this keeps the input fixed, and at 0 temperature (deterministic sampling), becomes equivalent to $pass@1$. In long-horizon tasks, error compounds irrespective of temperature, as shown in Figure 13. Further, we find that even single-step error rates can grow as the output length increases. While this might seem similar to prior work showing degradation in long-context (Zhou et al., 2025a), using counterfactual experiments, we show this is rather due to the model self-conditioning on its own generated mistakes.

Limitations. As with any “synthetic” task used for a controlled study of LLM capabilities, there are a few limitations of our setup. It does not reflect complexities and sources of error arising in real agentic tasks with a large number of possible actions. In such settings, the number of actions and the accuracy of each action can both vary based on the plan, requiring more careful consideration. It would be interesting future work to study the self-conditioning effect when doing diverse actions instead of repeating the same ones. Our results are observations about pretrained LLMs, and not inherent properties of transformers, so they might change with task-specific finetuning. Improvement on our task is necessary, but not sufficient for long-horizon execution on real-world tasks. Finally, our current task accuracy metric does not account for self-correction. In tasks where mistakes are acceptable and easy to undo, self-correction is a promising direction to improve long-horizon execution.

In the Appendix, we design experiments that dig deeper into our setup. First, in Section D we note that one step in our task actually requires three operations—retrieval of the value, reading the current state, and adding to it. Individually, we find much better accuracies at long-horizon execution of each of these operations, but taken together, due to increased turn complexity, errors grow much faster. Second, in Section C we show that the horizon length of different models can vary significantly at different turn complexities. This emphasizes the importance of contextualizing any claims about the length of task a model can complete with the complexity within a turn. Finally, in Section A we do a preliminary investigation of possible fixes, including self-correction prompting at each turn, and a simple context management technique. Specifically, we find significant improvements from simply removing historic context, which reduces the probability of errors appearing in context for self-conditioning. However, this exploits the Markovian nature of our task and would not work if we added dependencies to arbitrary previous states, such as in dynamic programming (Beniamini et al., 2025).

Outlook. Generative models that maintain accuracy over long horizons will be essential for creating simulated environments (Bruce et al., 2024) to train open-ended agents (Raad et al., 2024). Scaling up the length of tasks a model can complete would be a major step towards realizing the true potential of general agents. By showing that long-horizon execution can be studied on simple tasks, we hope to inspire more research on this capability.

Acknowledgements

We thank Maksym Andriushchenko, Nikhil Chandak, Paras Chopra, Dulhan Jayalath, Abhinav Menon, Sumeet Motwani, Ameya Prabhu, and Shashwat Singh for helpful feedback. AA was funded by the CHIPS Joint Undertaking (JU) under grant agreement No. 101140087 (SMARTY), and by the German Federal Ministry of Education and Research (BMBF) under the sub-project with the funding number 16MEE0444. AA thanks the International Max Planck Research School for Intelligent Systems (IMPRS-IS) and the European Laboratory for Learning and Intelligent Systems (ELLIS) PhD program for support. The authors gratefully acknowledge compute time on the Artificial Intelligence Software Academy (AISA) cluster funded by the Ministry of Science, Research and Arts of Baden-Württemberg.

Author Contributions

SG conceived the project. AS led the execution of the experiments with the help of AA, while SG led their planning with the help of AA, AS, and JG. SG and AA wrote the paper, while AS worked on the figures. JG and SS advised the project, providing valuable feedback throughout.

References

- Zeyuan Allen-Zhu. ICML 2024 Tutorial: Physics of Language Models, July 2024. Project page: <https://physics.allen-zhu.com/>.
- Anthropic. System card: Claude opus 4 & claude sonnet 4, May 2025. URL <https://www.anthropic.com/claude-4-system-card>. Covers Claude Sonnet 4 and Opus 4.
- Simran Arora, Sabri Eyuboglu, Aman Timalsina, Isys Johnson, Michael Poli, James Zou, Atri Rudra, and Christopher Re. Zoology: Measuring and Improving Recall in Efficient Language Models. In *The Twelfth International Conference on Learning Representations*, October 2023. URL <https://openreview.net/forum?id=LY3ukUANko>.
- Axel Backlund and Lukas Petersson. Vending-Bench: A Benchmark for Long-Term Coherence of Autonomous Agents. *arxiv:2502.15840[cs]*, February 2025. doi: 10.48550/arXiv.2502.15840. URL <http://arxiv.org/abs/2502.15840>.
- Jonas Becker, Lars Benedikt Kaesberg, Andreas Stephan, Jan Philip Wahle, Terry Ruas, and Bela Gipp. Stay Focused: Problem Drift in Multi-Agent Debate. *arxiv:2502.19559[cs]*, May 2025. doi: 10.48550/arXiv.2502.19559. URL <http://arxiv.org/abs/2502.19559>.
- Gal Beniamini, Yuval Dor, Alon Vinnikov, Shir Granot Peled, Or Weinstein, Or Sharir, Noam Wies, Tomer Nussbaum, Ido Ben Shaul, Tomer Zekharya, Yoav Levine, Shai Shalev-Shwartz, and Amnon Shashua. Formulaone: Measuring the depth of algorithmic reasoning beyond competitive programming, 2025. URL <https://arxiv.org/abs/2507.13337>.
- Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024.
- Ziyang Cai, Nayoung Lee, Avi Schwarzschild, Samet Oymak, and Dimitris Papailiopoulos. Extrapolation by Association: Length Generalization Transfer in Transformers. *arxiv:2506.09251[cs]*, August 2025. doi: 10.48550/arXiv.2506.09251. URL <http://arxiv.org/abs/2506.09251>.
- Siwei Chen, Anxing Xiao, and David Hsu. LLM-State: Open World State Representation for Long-horizon Task Planning with Large Language Model. *arxiv:2311.17406[cs]*, April 2024a. doi: 10.48550/arXiv.2311.17406. URL <http://arxiv.org/abs/2311.17406>.
- Yanan Chen, Ali Pesaranghader, Tanmana Sadhu, and Dong Hoon Yi. Can We Rely on LLM Agents to Draft Long-Horizon Plans? Let’s Take TravelPlanner as an Example. *arxiv:2408.06318[cs]*, August 2024b. doi: 10.48550/arXiv.2408.06318. URL <http://arxiv.org/abs/2408.06318>.
- Jingde Cheng. Why cannot large language models ever make true correct reasoning?, 2025. URL <https://arxiv.org/abs/2508.10265>.
- Junhyuk Choi, Yeseon Hong, Minju Kim, and Bugeun Kim. Examining identity drift in conversations of llm agents. *arXiv preprint arXiv:2412.00804*, 2024.

- Francois Chollet, Mike Knoop, Gregory Kamradt, and Bryan Landers. Arc prize 2024: Technical report. *arXiv preprint arXiv:2412.04604*, 2024.
- DeepSeek-AI, Aixin Liu, Bei Feng, et al. Deepseek-v3 technical report, 2025. URL <https://arxiv.org/abs/2412.19437>.
- Fabrizio Dell’Acqua, Edward McFowland III, Ethan R. Mollick, Hila Lifshitz-Assaf, Katherine C. Kellogg, Saran Rajendran, Lisa Kraye, François Candelon, and Karim R. Lakhani. Navigating the jagged technological frontier: Field experimental evidence of the effects of AI on knowledge worker productivity and quality. Working paper, Harvard Business School Technology & Operations Management Unit, 2023. URL <https://ssrn.com/abstract=4573321>. Also circulated as The Wharton School Research Paper; last revised 2023-09-27.
- Ying Fan, Yilun Du, Kannan Ramchandran, and Kangwook Lee. Looped Transformers for Length Generalization. In *The Thirteenth International Conference on Learning Representations*, October 2024. URL <https://openreview.net/forum?id=2edigk8yoU>.
- Richard E. Fikes and Nils J. Nilsson. Strips: A new approach to the application of theorem proving to problem solving. *Artificial Intelligence*, 2(3):189–208, December 1971. ISSN 0004-3702. doi: 10.1016/0004-3702(71)90010-5. URL <https://www.sciencedirect.com/science/article/pii/0004370271900105>.
- Gemini Team. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. Technical report, Google DeepMind, June 2025. URL https://storage.googleapis.com/deepmind-media/gemini/gemini_v2_5_report.pdf.
- Gemma-Team, Aishwarya Kamath, Johan Ferret, et al. Gemma 3 technical report, 2025. URL <https://arxiv.org/abs/2503.19786>.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- Jacob Hilton, Jie Tang, and John Schulman. Scaling laws for single-agent reinforcement learning. *arXiv preprint arXiv:2301.13442*, 2023.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekish, Fei Jia, and Boris Ginsburg. RULER: What’s the Real Context Size of Your Long-Context Language Models? In *First Conference on Language Modeling*, August 2024. URL <https://openreview.net/forum?id=kIoBbc76Sy>.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant Bhambri, Lucas Saldyt, and Anil Murthy. Llm’s can’t plan, but can help planning in llm-modulo frameworks. *arXiv preprint arXiv:2402.01817*, 2024.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020. URL <https://arxiv.org/abs/2001.08361>.
- Sheraz Khan, Subha Madhavan, and Kannan Natarajan. A comment on "the illusion of thinking": Reframing the reasoning cliff as an agentic gap. *arXiv preprint arXiv:2506.18957*, 2025.
- Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv preprint arXiv:2507.20534*, 2025.
- Kimi-Team, Yifan Bai, Yiping Bao, et al. Kimi k2: Open agentic intelligence, 2025. URL <https://arxiv.org/abs/2507.20534>.

- Thomas Kwa, Ben West, Joel Becker, Amy Deng, Katharyn Garcia, Max Hasin, Sami Jawhar, Megan Kinniment, Nate Rush, Sydney Von Arx, et al. Measuring ai ability to complete long tasks. *arXiv preprint arXiv:2503.14499*, 2025.
- Yann LeCun. Do large language models need sensory grounding for meaning and understanding? Slide deck, NYU Philosophy of Deep Learning debate, March 2023. URL https://drive.google.com/file/d/1BU5bV3X5w65DwSMapKcsr0ZvrMRU_Nbi/view. Includes slide “Autoregressive LLMs are Doomed.”.
- Nayoung Lee, Ziyang Cai, Avi Schwarzschild, Kangwook Lee, and Dimitris Papailiopoulos. Self-improving transformers overcome easy-to-hard and length generalization challenges. *arXiv preprint arXiv:2502.01612*, 2025.
- Derrick H Lehmer. Teaching combinatorial tricks to a computer. In *Proceedings of Symposia in Applied Mathematics*, pages 179–193. American Mathematical Society, 1960.
- Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as Algorithms: Generalization and Stability in In-context Learning. In *Proceedings of the 40th International Conference on Machine Learning*, pages 19565–19594. PMLR, July 2023. URL <https://proceedings.mlr.press/v202/li231.html>.
- Jack Lindsey, Wes Gurnee, Emmanuel Ameisen, Brian Chen, Adam Pearce, Nicholas L. Turner, Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, Jonathan Marcus, Michael Sklar, Adly Templeton, Trenton Bricken, Callum McDougall, Hoagy Cunningham, Thomas Henighan, Adam Jermy, Andy Jones, Andrew Persic, Zhenyi Qi, T. Ben Thompson, Sam Zimmerman, Kelley Rivoire, Thomas Conerly, Chris Olah, and Joshua Batson. On the biology of a large language model. *Transformer Circuits Thread*, 2025. URL <https://transformer-circuits.pub/2025/attribution-graphs/biology.html>.
- William Merrill and Ashish Sabharwal. The expressive power of transformers with chain of thought. *arXiv preprint arXiv:2310.07923*, 2023.
- METR. Measuring ai ability to complete long tasks, March 2025. URL <https://metr.org/blog/2025-03-19-measuring-ai-ability-to-complete-long-tasks/>.
- Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. *arXiv preprint arXiv:2410.05229*, 2024.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task generalization via natural language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*, 2021.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. In-context Learning and Induction Heads. *CoRR*, January 2022. URL <https://openreview.net/forum?id=nJ10GgImU0>.
- OpenAI. Gpt-5 system card, August 2025. URL <https://cdn.openai.com/gpt-5-system-card.pdf>. Canonical system card PDF.
- Shishir G Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and Joseph E Gonzalez. The berkeley function calling leaderboard (bfcl): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*, 2025.
- Michael Poli, Armin W. Thomas, Eric Nguyen, Pragaash Ponnusamy, Björn Deiseroth, Kristian Kersting, Taiji Suzuki, Brian Hie, Stefano Ermon, Christopher Re, Ce Zhang, and Stefano Massaro. Mechanistic Design and Scaling of Hybrid Architectures. In *Forty-First International Conference on Machine Learning*, June 2024. URL <https://openreview.net/forum?id=GDp7Gyd9nf>.
- Maria Abi Raad, Arun Ahuja, Catarina Barros, Frederic Besse, Andrew Bolt, Adrian Bolton, Bethanie Brownfield, Gavin Buttimore, Max Cant, Sarah Chakera, et al. Scaling instructable agents across many simulated worlds. *arXiv preprint arXiv:2404.10179*, 2024.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, 2023.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

- Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng Li, and Yueting Zhuang. TaskBench: Benchmarking large language models for task automation. In *Proceedings of the 38th International Conference on Neural Information Processing Systems*, volume 37 of *NIPS '24*, pages 4540–4574, Red Hook, NY, USA, June 2025. Curran Associates Inc. ISBN 979-8-3313-1438-5.
- Parshin Shojadeh, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity, 2025. URL <https://arxiv.org/abs/2506.06941>.
- David Silver and Richard S Sutton. Welcome to the era of experience. *Google AI*, 1, 2025.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2021.
- Simeng Sun, Cheng-Ping Hsieh, Faisal Ladhak, Erik Arakelyan, Santiago Akle Serano, and Boris Ginsburg. L0-Reasoning Bench: Evaluating Procedural Correctness in Language Models via Simple Program Execution. *arxiv:2503.22832[cs]*, April 2025. doi: 10.48550/arXiv.2503.22832. URL <http://arxiv.org/abs/2503.22832>.
- Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. Long Range Arena : A Benchmark for Efficient Transformers. In *International Conference on Learning Representations*, October 2020. URL <https://openreview.net/forum?id=qVyeW-grC2k>.
- Karthik Valmeekam, Kaya Stechly, Atharva Gundawar, and Subbarao Kambhampati. A Systematic Evaluation of the Planning and Scheduling Abilities of the Reasoning Model o1. *Transactions on Machine Learning Research*, December 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=FkKBxp0FhR>.
- Joshua Vendrow, Edward Vendrow, Sara Beery, and Aleksander Madry. Do large language model benchmarks test reliability? *arXiv preprint arXiv:2502.03461*, 2025.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In *International Conference on Machine Learning*, pages 11080–11090. PMLR, 2021.
- xAI. Grok 4 model card, August 2025. URL <https://data.x.ai/2025-08-20-grok-4-model-card.pdf>.
- Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. Travelplanner: A benchmark for real-world planning with language agents, 2024. URL <https://arxiv.org/abs/2402.01622>.
- An Yang, Anfeng Li, Baosong Yang, et al. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.
- Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. tau-bench: A benchmark for tool-agent-user interaction in real-world domains. *arXiv preprint arXiv:2406.12045*, 2024.
- Haoyu Zhou, Mingyu Ding, Weikun Peng, Masayoshi Tomizuka, Lin Shao, and Chuang Gan. Generalizable Long-Horizon Manipulations with Large Language Models. *arxiv:2310.02264[cs]*, October 2023. doi: 10.48550/arXiv.2310.02264. URL <http://arxiv.org/abs/2310.02264>.
- Yang Zhou, Hongyi Liu, Zhuoming Chen, Yuandong Tian, and Beidi Chen. Gsm-infinite: How do your llms behave over infinitely increasing context length and reasoning complexity? *arXiv preprint arXiv:2502.05252*, 2025a.
- Yang Zhou, Hongyi Liu, Zhuoming Chen, Yuandong Tian, and Beidi Chen. GSM- ∞ : How Do your LLMs Behave over Infinitely Increasing Reasoning Complexity and Context Length? In *Forty-Second International Conference on Machine Learning*, June 2025b. URL <https://openreview.net/forum?id=n52yyvEwPa>.

Minjun Zhu, Qiujie Xie, Yixuan Weng, Jian Wu, Zhen Lin, Linyi Yang, and Yue Zhang. AI Scientists Fail Without Strong Implementation Capability. *arxiv:2506.01372[cs]*, June 2025. doi: 10.48550/arXiv.2506.01372. URL <http://arxiv.org/abs/2506.01372>.

Appendix

Contents

A Investigating proposed fixes for long-horizons tasks	17
A.1 Turn-wise Verification Prompting	17
A.2 Context Management	17
B Can parallel test-time compute scaling match thinking?	18
C Number of turns vs turn complexity	18
D Deconstructing error in retrieve-then-compose	19
E Experimental Setup	21
E.1 Task Details	21
E.2 Prompting	21
E.3 Prompting For Thinking Models	21
E.4 Model Specifications	21
E.5 Compute Details	22
F Format Following Failures	22
G Chain-of-Thought Self-Conditioning	23
H Proof and Analysis of Proposition 1	24
H.1 Implications for Horizon Length ($H_{0.5}$)	24

A Investigating proposed fixes for long-horizons tasks

A.1 Turn-wise Verification Prompting

We investigate whether self-simulation can be mitigated by explicitly prompting the model to perform active self-correction. At each turn, we instruct the model to first re-validate its previously reported state and, if required, recalculate the full historical sum before processing the current turn’s keys.

The results, shown in Figure 8, are mixed. For the Gemma3 family with CoT, this prompt provides an initial boost in accuracy, successfully breaking the self-simulation loop in early turns. However, the self-verification process significantly increases the number of tokens generated per turn, causing the model to exhaust its context window much sooner, which leads to a sharper performance collapse in later stages. In contrast, the Qwen3 thinking models show negligible improvement. Manually inspecting of their reasoning traces, we find that these models, likely due to their fine-tuning, overthink and frequently fail at the verification step itself, sometimes making arithmetic errors even during their re-calculation process.

These findings suggest that prompting self-correction may not be a viable solution. It is computationally expensive, incurring a context-length penalty, and is itself a complex, error-prone execution task that models may not be able to perform reliably.

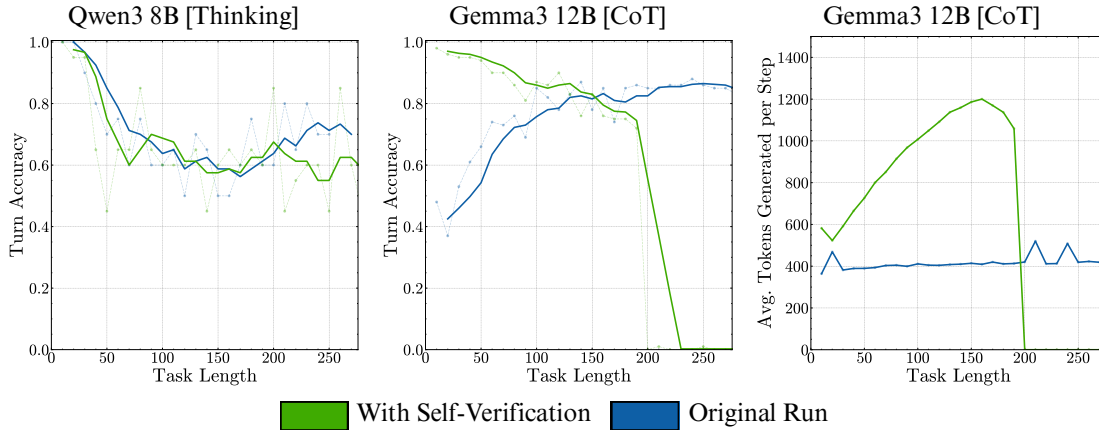


Figure 8: Self-verification prompting. Prompting to self-verify does not suffice to fix the self-conditioning effect completely. It leads to overthinking in thinking models and increases the amount of tokens required per turn, leading to faster context consumption in CoT models.

A.2 Context Management

Another natural mitigation strategy is to limit the model’s exposure to its own past errors in its history. We operationalize this using a simple sliding context window, which is particularly well-suited for Markovian tasks like ours. This approach maintains only the N most recent turns in the model’s context. The rationale is that a smaller context window reduces the probability of the model observing a lot of its own past failures, thereby breaking the negative feedback loop of self-conditioning.

As shown in Figure 9(a), performance improves significantly as the context window size is reduced, allowing models to sustain execution for longer horizons. While a fixed sliding window is only applicable to tasks without long-range dependencies, this result validates a more general principle: active context management designed to minimize the accumulation of errors in the context is a promising direction for improving long-horizon reliability in LLM agents.

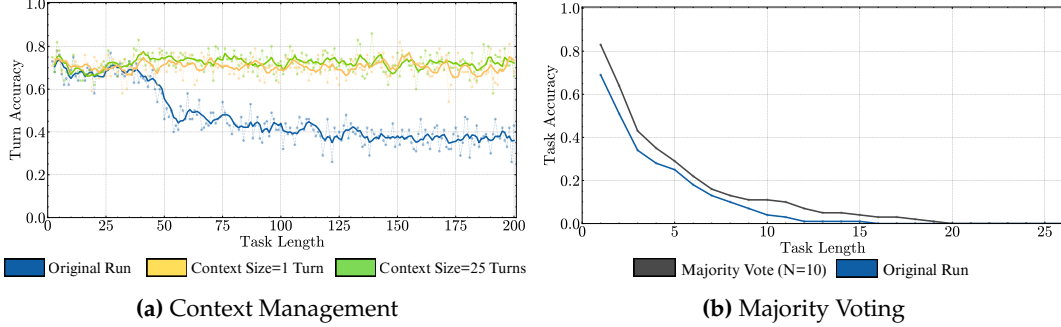


Figure 9: Context Management and Majority Voting on Gemma3 12B. Controlling context size reduces the self-conditioning effect, but relies on the Markovian nature of our task. Majority voting at $K=1$ provides only minimal improvements over the baseline.

B Can parallel test-time compute scaling match thinking?

We also experiment to validate if parallel scaling in test-time compute can achieve the same improvements as thinking. We verify this by testing if parallel majority voting can replicate the gains from either model scale or sequential computation (*thinking*). To create a fair comparison, we sample multiple outputs from a non-thinking Gemma3 model at each turn, with the number of samples set to match the average token count of its CoT counterpart. The final answer is determined by a majority vote over these parallel generations. From the results in Figure 10 and 9(b), we see that while majority voting yields a marginal performance improvement, it is insufficient to match the reliability of a larger, non-thinking model, let alone the substantial gains from using CoT reasoning. This suggests that for long-horizon execution, sequential computation provides an advantage that parallel test time scaling cannot match. This contrasts with findings in other domains, such as math or common-sense reasoning, where parallel sampling with self-consistency has been shown to be highly competitive (Snell et al., 2024).

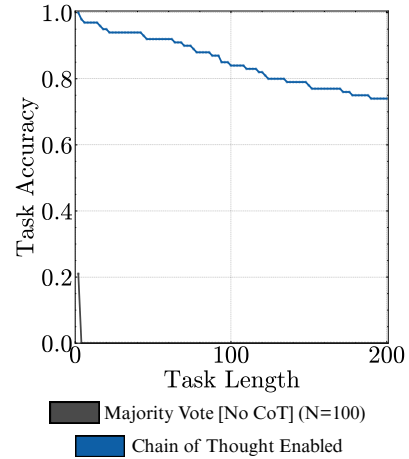
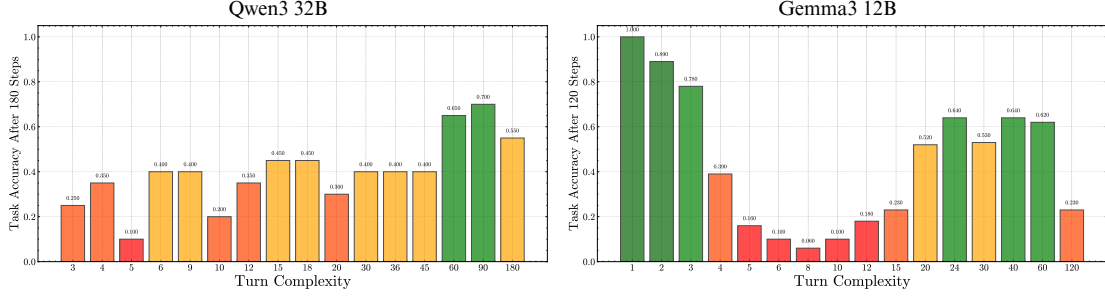


Figure 10: Parallel test time scaling on Gemma3 12B at $K=2$. Majority voting with the same amount of tokens as CoT traces does not nearly match the performance with CoT.

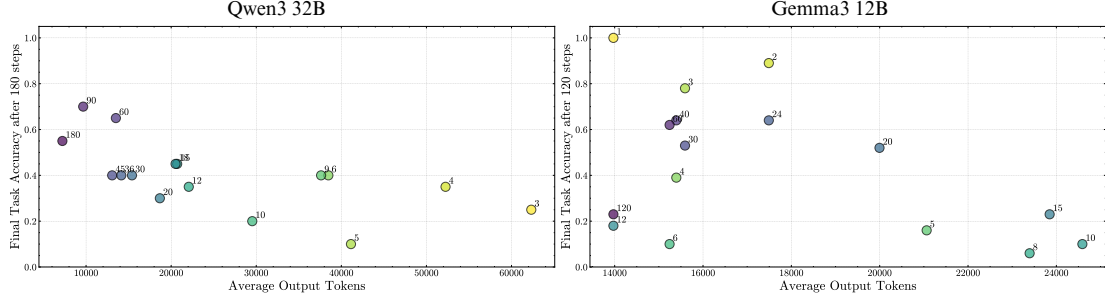
C Number of turns vs turn complexity

In our experiments, we show that we can increase the length of the task needed to be performed by either (1) increasing the number of turns or (2) increasing the turn complexity, i.e., providing more inputs in the same turn. To investigate the relationship between these two axes, we perform an experiment where a model has to perform a fixed number of operations while varying the turn complexity. A higher turn complexity means the model requires fewer turns to reach the fixed number of operations. Results in Figure 11 indicate there is no strict turn complexity that is consistently the best across model families. Rather, we found that different models behaved quite differently for the same turn complexities. Qwen3 32B seems to show poorer performance at lower turn complexities, indicating that it is unable to perform well over a large number of turns, even if the turns are simple themselves. Gemma3 12B shows a different trend. It reaches accuracy peaks at either extreme of the turn complexity spectrum, failing badly at mid-level turn complexities. This indicates it suffers when the turn complexity and the number of turns are both sufficiently high.

Another axis of evaluating the number of turns vs turn complexity trade-off is the test-time compute used. From an economic view, increasing the number of turns increases the overall



(a) For the same number of total steps, different turn complexities lead to different outcomes. We find no trend across families.



(b) Average output tokens used to complete the execution vs final accuracy. We see that for Qwen3 32B, more turns lead to more token usage, even at lower turn complexities, pointing to overthinking. Gemma3 12B, on the other hand, uses less tokens for very low turn complexity or very high turn complexity.

Figure 11: Relation between the turn complexity and the number of turns.

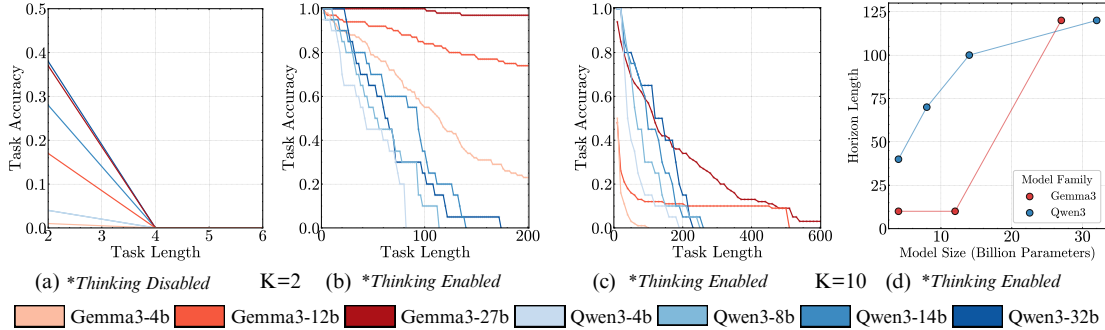


Figure 12: Scaling trends hold even enabling Sequential Test Time compute. We compare model performance with thinking disabled (a) against thinking enabled (b, c) at varying turn complexities. (a) Without thinking, all models fail to execute even two steps ($K = 2$) in a single turn. (b) In contrast, enabling thinking prevents this performance collapse, with all models successfully handling $K = 2$. (c) When the turn complexity is further increased to $K = 10$, performance degrades, but a clear scaling trend emerges. (d) This trend is explicitly shown, illustrating that for complex turns, the horizon length increases consistently with model size, reinforcing the benefits of scaling model size even when thinking is enabled.

cost of inference. We can lower the number of turns by increasing the turn complexity, but that would result in an increase in the per-turn inference cost, as a result of the added complexity. For the same experiment above, we track the number of output tokens used for computation (including thinking tokens) per sample, and again find diverging results for each family.

D Deconstructing error in retrieve-then-compose

To further isolate the source of execution error, we decompose our task into its two constituent operations—retrieval and addition—and evaluate models on them individually:

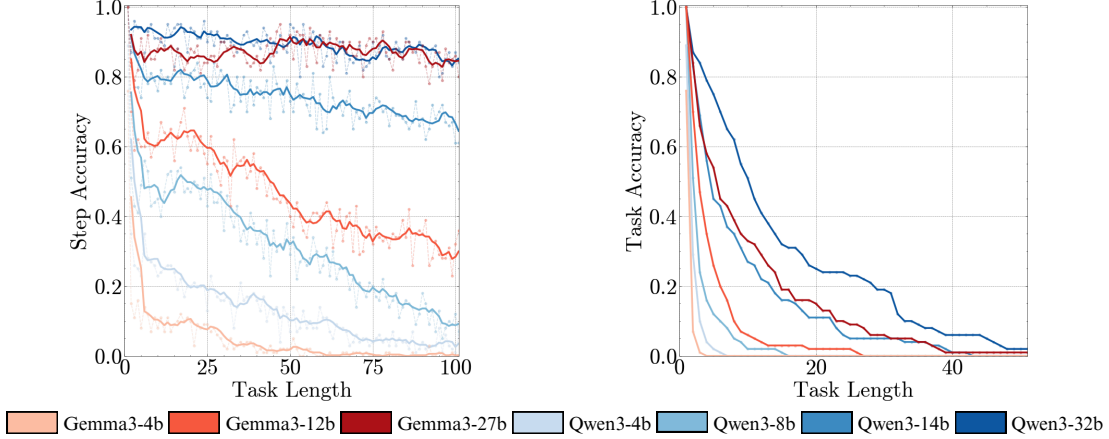


Figure 13: Temperature does not impact the trends observed. We reproduce the same trends in Figure 4, when running with temperature 0.

- **Retrieval-Only Task.** A stateless task where, at each turn, the model is given a key and must simply return the corresponding integer value from the dictionary. No running sum is maintained.
- **Addition-Only Task.** A stateless task where, at each turn, the model is given two random integers to add. No running sum is maintained. This isolates the arithmetic component.
- **Prefix-sum Task.** A stateful task where, at each turn, the model is given an integer directly and must add it to its previously reported running sum. This isolates the arithmetic and state-tracking component.

From Section D, we can observe that models achieve near-perfect performance on the stateless retrieval and addition task, indicating that neither simple dictionary lookup nor addition is a significant source of error. In contrast, the prefix sum task, while significantly better than our task, still exhibits a slow degradation over time.

This leads to two key insights. First, the difficulty lies not in the atomic operations themselves, which models perform with high accuracy in isolation over long horizons. Second, this suggests that the primary source of degradation is the **state-management** component of the task. While stateless retrieval and addition are trivial, the requirement to reliably maintain and update a running sum introduces higher chances of error. This suggests that the models struggle with the requirement to concurrently manage information lookup and state updates.

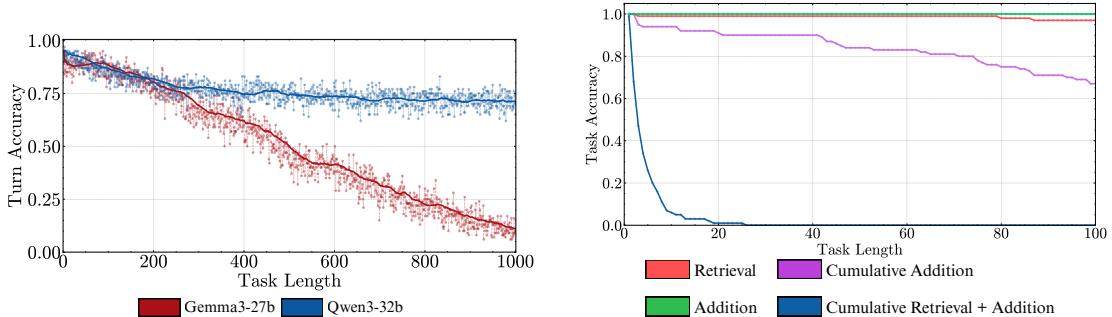


Figure 14: Analysis of execution failures. (a) Self-conditioning effect emerges as tasks get longer. Even for models that ace the task at a task length of 100, the Turn Accuracy drops constantly as we further increase the turns. (b) Models are good at the tasks individually, but not on their composition. State tracking introduces additional difficulty.

E Experimental Setup

E.1 Task Details

We create a dictionary where keys consist of common five-letter English words, and the values consist of integers uniformly sampled from -99 to 99 . The range of values is deliberately kept to be large to minimize the chance of an assistant being wrong in an earlier turn correcting its response by pure accident. For our experiment, we first create a fixed set of 100 keys. Then, we create multiple rollouts (samples) of 50,000 steps. For each rollout, we uniformly sample a separate set of values to be assigned to each key, in order to increase experimental breadth. Next, at each step, we uniformly sample a key to be provided at that step with replacement. This gives us a list of keys to be processed in order, which is exactly the plan to be executed by the agent.

To account for the turn complexity (K), we group K consecutive keys together and represent them as *one turn*. Thus, we can fully specify our intended evaluation by (1) specifying the number of samples (rollouts) needed, (2) the turn complexity, K , and (3) the number of turns required. This gives us a superset of data from which we sample rollouts to use in our experiments. For the Qwen3 and Gemma3 families, we sample 100 rollouts. For frontier models, due to cost limitations, we sample 20-50 rollouts. To ensure consistency in evaluation, we provide the same rollouts to each model.

E.2 Prompting

Each LLM is provided a standardized prompt describing the task at the start of the conversation. This prompt specifies the dictionary containing the five-letter word keys and their corresponding values. Further, the prompt specifies the number of keys that will be provided to the LLM at each subsequent turn. To ensure format following, the prompt also contains few-shot examples of different scenarios. Finally, the LLM is required to provide the running sum after each turn in `<answer>` tags. An example conversation is shown below.

E.3 Prompting For Thinking Models

To enable models to use chain-of-thought prompting, we add the line “Think step by step before answering.” to the prompt and added CoT traces to the in-context examples. We found that models stop performing CoT reasoning after few turns, as it starts conditioning on the answer format in its history. Thus, we end up *including the chain-of-thought trace* in the conversation history, to ensure the model does not forget the CoT instruction. This is a trade-off we had to make as it increases the input context of the LLM, however, it was essential to ensure instruction following. Thinking models provided their reasoning in `<think>` tags, which were removed from the conversation history. No other changes were needed to make the thinking models follow instructions.

E.4 Model Specifications

For chain-of-thought prompting, we set the **per-turn** output token limit to 10,000 tokens, and for thinking models, the token limit is set to 32,000 tokens, consistent with token limits provided by OpenRouter. We ensure that these token limits are sufficient to complete the required computations.

We use a temperature of 0.6 and a top-p value of 0.95 for all Gemma models. For Qwen, we use a temperature of 0.6 and a top-p value of 0.95 for thinking mode and a temperature of 0.7 and a top-p value of 0.8 for non-thinking as recommended in their documentation.¹

¹<https://huggingface.co/Qwen/Qwen3-32B>

E.5 Compute Details

All experiments were conducted on machines equipped with 4x NVIDIA A100 GPUs with 40/80GB of memory. Frontier model evaluations were performed using OpenRouter.

Starting Prompt:

You are an AI assistant. I will provide you with a dictionary and then give you keys in groups of 2. Your task is to keep a running total (starting from 0) by adding the values associated with the keys I provide.

In each turn, I'll provide 2 keys (comma-separated). Respond with the current running sum, enclosed in <answer> tags.

Examples:

Dictionary to maintain: 'apple': 5, 'banana': 0, 'cherry': 7, 'grape': -4, 'kiwi': 2, 'mango': -1

Example 1: keys in groups of 2

User: apple, banana

Assistant: <answer>5</answer>

User: cherry, grape

Assistant: <answer>8</answer>

User: kiwi, mango

Assistant: <answer>9</answer>

Example 2: keys in groups of 3

User: apple, banana, cherry

Assistant: <answer>12</answer>

User: grape, kiwi, mango

Assistant: <answer>9</answer>

Example 3: keys in groups of 6

User: apple, banana, cherry, grape, kiwi, mango

Assistant: <answer>9</answer>

Now, here is the actual task:

Dictionary to maintain:

'doubt': -64, 'alone': 46, 'adult': 84, 'fault': -19, 'brain': -45, 'blind': 68, ...
'coach': -31, 'alarm': 88, 'could': 25, 'cable': -32

Ready to start!

IMPORTANT: DO NOT OUTPUT ANY OTHER TEXT OUTSIDE ANSWER TAGS. Only provide the final running sum OF ALL TURNS in <answer> tags.

User: alarm, coach

Assistant: <answer>57</answer>

User: doubt, cable

Assistant: <answer>-39</answer>

F Format Following Failures

In any LLM evaluation, format following failures are a common source of error that is often neglected. In our experiments, any model can have 2 types of format following failures: (1) They do not provide <answer> tags in their answer, and (2) They do not provide a valid integer within <answer> tags. We take multiple steps to minimize format following failures. We ensure clarity in the starting prompt with clear format instructions, as well as few-shot examples. To empirically verify that model errors on our task are actually execution errors and not just format following errors in disguise, for each experiment, we also track the format failure fraction: the fraction of samples that do not correctly follow the format, with the failure being either (1) or (2). It is important to note that while we try to minimize any such error to the best of our abilities, we still count format following as a limitation of the model and hence a source of error.

Our results for format following failures are presented in Figure 15. For the experiment presented in Figure 4, we observe that smaller models are more susceptible to format failures, with the Qwen3 family in particular being worse at following format instructions. Overall, the fraction

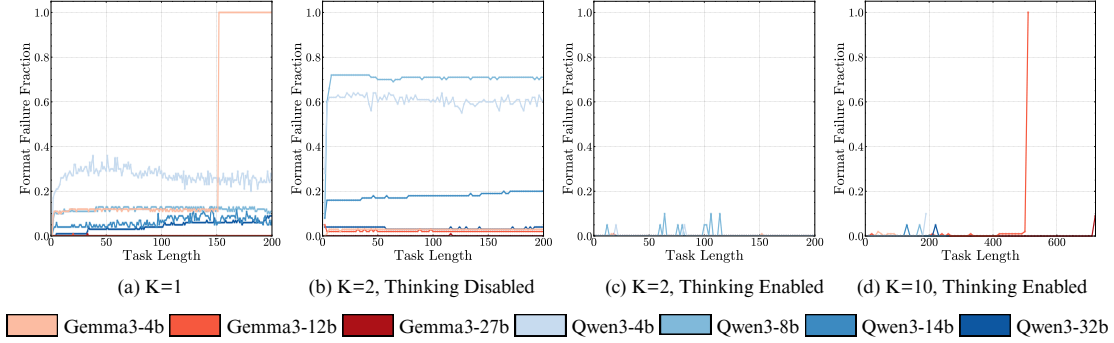


Figure 15: Analysis of format following failures. We analyze the fraction of errors attributed to incorrect format following for the experiments presented in Section 3. Overall, format adherence is high and not the primary source of execution errors.

for format following errors is low (around 0.1), with the Qwen3-8B being an exception. We find that the error here actually comes from the model trying to cheat, and do the entire summation *inside* the `<answer>` tags (For example, `<answer>39 + 51 = 90</answer>`). This is explicitly forbidden as we do not allow chain-of-thought or thinking in this experiment, and thus we count this as an error. Gemma3 4B fails at later turns due to a full context, however that does not affect any results, as its accuracies drop much earlier.

For the experiments presented in Figure 15, we find the Qwen3 family to be prone to format following errors in the case where we have thinking disabled for $K = 2$. We again find this to be the consequence of models trying to cheat and use extra tokens for computation inside the answer tags. This is fixed by enabling thinking. Following this, the errors in format following become negligible. At $K = 10$, we see Gemma3 12B sharply rise to a format failure fraction of 1.0, again due to a full context window.

G Chain-of-Thought Self-Conditioning

While our self-conditioning analysis provides clear insights for thinking models, extending this to models using Chain-of-Thought (CoT) presents some significant methodological challenges.

First, a fundamental prerequisite for reliable CoT reasoning is the inclusion of prior CoT traces in the context history. As we observed with the Gemma3 models, they often condition on the format of the context; if prior turns lack CoT traces, the models cease to generate them, even when explicitly instructed to do so. Consequently, this experiment for CoT must include the full reasoning trace for every preceding turn. This requirement immediately makes the setup practically infeasible, as the verbose nature of CoT traces would rapidly exhaust the context window limits of even frontier models.

Second, even if context length were not a constraint, the process of injecting controlled errors into CoT histories is not straightforward. A naive approach of only altering the final answer while preserving the original, correct CoT trace creates an unfaithful history. When conditioned on a history where reasoning and conclusions are contradictory, the model is no longer being tested on its execution reliability but on how it resolves inconsistency—it might learn to distrust its own reasoning, introducing a confounding variable.

The alternative is to programmatically generate flawed CoT traces. We implemented and experimented with this; however, this introduces its own complexities. For our simple task, there are multiple distinct points of failure within a single trace: an error in the retrieval step (looking up an incorrect value) or an error in the composition step (an arithmetic mistake). A controlled experiment would need to systematically manage the type, frequency, and location of these injected errors, making the setup intractable. Even establishing a “perfectly correct” (Induced Error Rate = 0) baseline history is problematic. A model might have a CoT trace with flawed reasoning (e.g., a minor calculation error that cancels out), which we then replace with the correct final answer. Such a history is also unfaithful.

Given these challenges—the practical infeasibility due to context length and the difficulty of designing a faithful error injection mechanism, we limit our self-conditioning analysis to non-thinking and thinking models.

H Proof and Analysis of Proposition 1

Proposition 1. *Assuming a constant per-step accuracy p and no self-correction, the horizon-length H at which a model can achieve a success rate s is given by:*

$$H_s(p) = \frac{\ln(s)}{\ln(p)}$$

Proof. Let p be the constant probability of successfully executing a single step. Under the assumption of no self-correction, a task of length H is successful only if all H independent steps are executed correctly. The probability of this joint event, $P(\text{success}, H)$, is the product of the individual step probabilities:

$$P(\text{success}, H) = \underbrace{p \times p \times \cdots \times p}_{H \text{ times}} = p^H$$

This is equivalent to the Task Accuracy at turn H , i.e., $\text{TA}(H) = p^H$. We define the horizon-length H as the number of turns at which the probability of success equals a desired rate s . Therefore, we set our expression for the success probability equal to s :

$$p^H = s$$

Solving for H ,

$$\begin{aligned} \ln(p^H) &= \ln(s) \Rightarrow H \cdot \ln(p) = \ln(s) \\ H_s(p) &= \left\lceil \frac{\ln(s)}{\ln(p)} \right\rceil \approx \frac{\ln(s)}{\ln p} \end{aligned}$$

This completes the proof. □

H.1 Implications for Horizon Length ($H_{0.5}$)

We can apply this general result to our specific metric, the Effective Task Length ($H_{0.5}$), which is defined as the number of turns at which Task Accuracy drops to $s = 0.5$,

$$H_{0.5}(p) = \left\lceil \frac{\ln(0.5)}{\ln p} \right\rceil = \left\lceil -\frac{\ln(2)}{\ln p} \right\rceil$$

For analysis, we use the continuous approximation:

$$H_{0.5}(p) \approx -\frac{\ln(2)}{\ln p}$$

Sensitivity to Small Changes in Step Accuracy. This formulation allows us to analyze the sensitivity of the horizon length to small improvements in per-step accuracy by taking the derivative with respect to p ,

$$\frac{dH_{0.5}}{dp} = -\ln(2) \cdot \left(-\frac{1}{(\ln p)^2} \cdot \frac{1}{p} \right) = \frac{\ln 2}{p(\ln p)^2}$$

This implies that a small change in accuracy Δp results in a change in horizon length $\Delta H_{0.5}$ of,

$$\Delta H_{0.5} \approx \frac{\ln 2}{p(\ln p)^2} \Delta p$$

Near-Perfect Accuracy Regime. The effect is most dramatic when accuracy is already high. For near-perfect accuracy, let $p = 1 - \varepsilon$ where $\varepsilon \ll 1$. Using the Taylor approximation $\ln(1 - \varepsilon) \approx -\varepsilon$, we can simplify the expression for $H_{0.5}$,

$$H_{0.5} \approx -\frac{\ln(2)}{\ln(1 - \varepsilon)} \approx -\frac{\ln(2)}{-\varepsilon} = \frac{\ln 2}{\varepsilon} = \frac{\ln 2}{1 - p}$$

The sensitivity in this regime becomes,

$$\frac{dH_{0.5}}{dp} \approx \frac{\ln 2}{(1 - p)^2} \Rightarrow \Delta H_{0.5} \approx \frac{\ln 2}{(1 - p)^2} \Delta p$$

This demonstrates that as $p \rightarrow 1$, the improvement in horizon length for a fixed gain in step accuracy grows quadratically, highlighting the compounding benefits of scale.