

COMPUTATIONAL REASONING OF LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

With the rapid development and widespread application of Large Language Models (LLMs), multidimensional evaluation has become increasingly critical. However, current evaluations are often domain-specific and overly complex, limiting their effectiveness as cross-domain proxies for core capabilities. To address these limitations and enable a unified and simple evaluation framework, an ideal proxy task should target a basic capability that generalizes across tasks and is independent of domain-specific knowledge. Turing machine provides a powerful theoretical lens by reducing complex processes to basic, domain-agnostic computational operations. This perspective offers a principled framework for evaluating foundational computational abilities essential to a wide range of tasks, particularly those involving complex, multi-step reasoning such as mathematics. Motivated by this abstraction, we introduce **Turing Machine Bench**, a benchmark designed to assess the ability of LLMs to **strictly follow rules** and **accurately manage internal states** for multi-step, referred to as **computational reasoning**. TMBench incorporates four key features: self-contained and knowledge-agnostic reasoning, a minimalistic multi-step structure, controllable difficulty, and a solid theoretical foundation based on Turing machine. Empirical results demonstrate that TMBench serves as an effective proxy for evaluating computational reasoning on representative LLMs. It produces clear step-wise accuracy curves, revealing LLMs’ ability to execute multi-step reasoning processes. By analyzing performance trends across TMBench and established reasoning benchmarks, we find strong correlations with real-world tasks, bridging real-task evaluation with basic ability assessment. These findings suggest that TMBench holds potential as a cross-domain dimension for evaluating reasoning in LLMs. Code and data are available at Repo.

1 INTRODUCTION

Recent progress in pre-training, post-training, and scaling has significantly advanced the capabilities of LLMs. These models now demonstrate remarkable capabilities not only in traditional natural language processing tasks, such as text classification and machine translation Devlin et al. (2019), but also in more complex cognitive domains, including advanced reasoning Rein et al. (2024); Huang et al. (2024), code generation Jimenez et al. (2023); Chen et al. (2021a), instruction following La Malfa et al. (2024), multimodal understanding Wu et al. (2023), and even scientific discovery Boiko et al. (2023). While these advancements significantly broaden the scope of potential applications, the growing complexity and diversity of cognitive tasks undertaken by LLMs present substantial challenges in accurately and reliably evaluating their true capabilities. The intelligence of LLMs can be mainly evaluated along two key dimensions: (1) knowledge and comprehension, and (2) reasoning and decision-making. First, knowledge and comprehension benchmarks primarily focus on the ability of LLM to apply its internal knowledge to understand and solve user questions. These skills are typically evaluated through large-scale question answering tasks (e.g., MMLU Hendrycks et al. (2020), TriviaQA Joshi et al. (2017)). Second, reasoning is the process of inferring and justifying conclusions from a set of premises Arkoudas (2023). In particular, reasoning encompasses a broad range of tasks, including logical deduction (e.g., Big-Bench Hard bench authors (2023)), mathematical reasoning (e.g., AIME MAA (2024)), formal proof (e.g., miniF2F Zheng et al. (2021), ProofNet Azerbayev et al. (2023)), programming (e.g., HumanEval Chen et al. (2021a)) and computational reasoning (e.g., Arithmetic Lai et al. (2024), CodeSimulation La Malfa et al. (2024)). Specifically, computational

reasoning refers to the ability to accurately interpret formal rules and execute multi-step computational operations without external tools. This ability is fundamental to modern sciences, which are generally founded on rule-based systems. Disciplines such as mathematics, physics, and chemistry rely on axioms, theorems, and laws, collectively referred to as rules. At the most fundamental level, arithmetic represents the simplest form of such rules. Although LLMs are not inherently proficient at direct numerical computation, they have demonstrated the ability to perform accurate arithmetic by simulating a Turing machine designed for arithmetic Lai et al. (2024). This suggests that LLMs have the potential to faithfully execute formal rules in multi-step reasoning tasks. By assessing the results across these two dimensions, we obtain a more comprehensive understanding of an LLM’s intelligence. However, these dimensions often overlap in practice. For instance, solving a mathematical problem may still rely on background knowledge encoded during pre-training, blurring the distinction between memorization (knowledge) and computation (reasoning).

To evaluate the computational reasoning ability of LLMs, we turn to the foundational paradigm of theoretical computer science: the Turing machine. Its minimalist design and computational universality provide a precise framework to measure a model’s ability to follow rules and deterministic state transitions. We introduce **TMBench**, a benchmark for evaluating computational reasoning of LLMs by simulating the operation of an m -Tag Turing machine. The m -Tag machine consists of production rules and a dynamic queue. At each step, it reads the queue head, appends the corresponding production symbols to the tail, and deletes the first m characters. Iterating this process and comparing the model’s execution trace with ground truth yields a quantitative measure of reasoning ability as an executor. TMBench has four key features: (1) **Self-contained reasoning**: all tasks are solvable from first principles without external knowledge; (2) **Minimalistic multi-step structure**: each step is necessary, interpretable, and yields a verifiable result; (3) **Controllable difficulty**: rule sets, m values, step counts, and input lengths can be tuned to vary complexity; (4) **Computational generality**: any computable function can be represented by a Turing machine, making TMBench a principled proxy for rule-following and state management. Our main contributions are summarized as follows:

1. We propose the TMBench to evaluate the computation reasoning ability of LLMs. Evaluations are conducted on a broad range of recent open-source models, from 0.6B to 671B parameters, as well as proprietary models such as Gemini, GPT, Grok, and Claude.
2. Empirical results demonstrate that TMBench strongly correlates with established reasoning benchmarks and real-world task performance, highlighting its practical relevance. It effectively captures LLMs’ multi-step reasoning abilities through clear step-wise accuracy trends, making it a promising cross-domain metric for evaluating computational reasoning.
3. To investigate TMBench’s characteristics, we conduct ablation studies on unbounded-step execution, decoding temperature, alphabet type, task difficulty and SFT. Results show that LLM inevitably fails with increasing steps due to its autoregressive nature. Its stable performance across different alphabets suggests reliance on reasoning rather than statistics. Varying deletion counts enables fine-grained control of task difficulty, demonstrating the benchmark’s scalability.

2 RELATED WORK

2.1 LLM BENCHMARKS

Benchmarks for large language models can be primarily categorized into two types: closed-ended benchmarks, which provide definitive answers, and open-ended benchmarks, which rely on human preference. For closed-ended benchmarks, a wide range of topics is addressed, including language understanding, mathematics, coding, reasoning, hallucination, toxicity, and stereotypes. Notable benchmarks in this category include MMLU Hendrycks et al. (2020), HellaSwag Zellers et al. (2019b), GSM-8K Cobbe et al. (2021a), HELM Liang et al. (2022), BigBench bench authors (2023), AGIEval Zhong et al. (2023), HumanEval Chen et al. (2021b), and ToxicChat Lin et al. (2023). Beyond closed-ended questions, benchmarks also include open-ended questions that with human preference. These questions are typically assessed through expert ratings or crowd-sourced evaluations. The recent trend includes utilizing GPT-4 for approximating human judgment Chiang & Lee (2023), with notable instances being MT-Bench Zheng et al. (2023) and AlpacaEval Li et al.

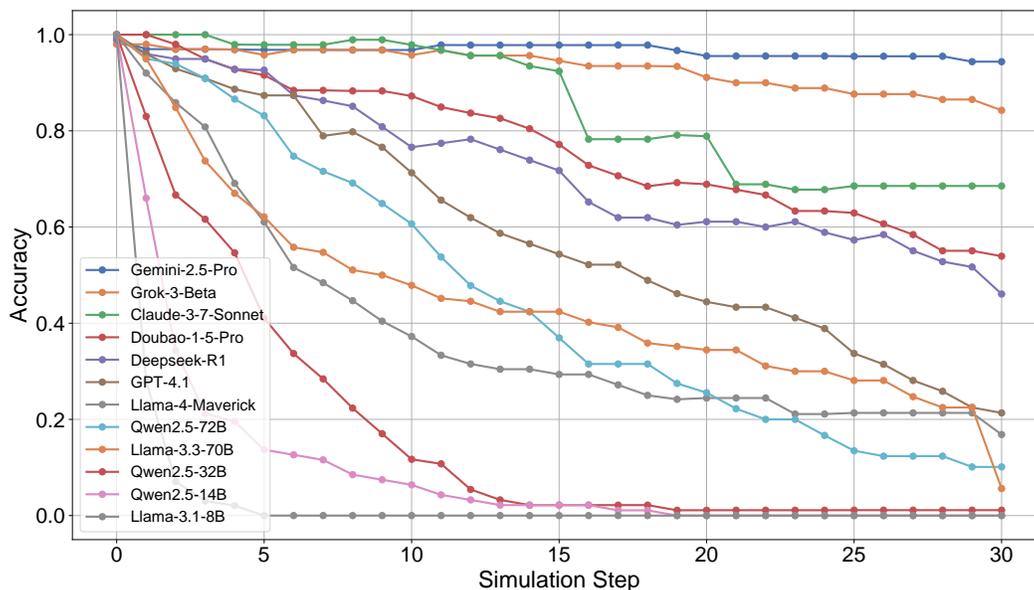


Figure 1: Illustration of the multi-step performance curve on TMBench across a diverse range of both open-source and proprietary LLMs. Proprietary models demonstrate advantages in computational reasoning abilities, but accuracy still decreases as steps increase.

(2023). In addition to static benchmarks, dynamic benchmarks featuring continuously updated questions are also available. These questions may be drawn from annual examinations, weekly online contests such as Codeforces Li et al. (2022); Huang et al. (2023), or collaborative platforms like ChatbotArena Zheng et al. (2023). Some studies have explored leveraging live human interaction for reinforcement learning from human preferences Bai et al. (2022); Ouyang et al. (2022); Touvron et al. (2023). In this paper, we introduce Turing Machine Bench, a self-contained, multi-step benchmark for computation reasoning that avoids the influence of external knowledge.

2.2 MULTI-STEP REASONING

Recent large language models demonstrate remarkable performance in solving complex reasoning tasks Lewkowycz et al. (2022). Multi-step reasoning, essential for tackling such tasks, has gained significant attention, with numerous methods proposed to enhance the reasoning capabilities of LLM. Chain-of-Thought (CoT) Wei et al. (2022b) and scratchpad Nye et al. (2021) encourages step-by-step reasoning by breaking down complex problems into manageable steps, while CoT’s extension, self-consistency Wang et al. (2022), improves performance by sampling multiple reasoning paths and selecting the most consistent solution through majority voting. These methods have been extensively evaluated across a variety of benchmarks, including commonsense reasoning Lourie et al. (2021); Geva et al. (2021), mathematical reasoning Cobbe et al. (2021b); Hendrycks et al. (2021), symbolic reasoning Luo et al. (2023); Han et al. (2022); Patel et al. (2024); Parmar et al. (2024), logical reasoning Liu et al. (2020); Tafjord et al. (2020), and multi-modal reasoning Zellers et al. (2019a); Xiao et al. (2021); Lu et al. (2022); Handa et al. (2024). These advancements have enhanced the diversity and reliability of multi-step reasoning evaluations. However, the varying levels of difficulty across steps in most reasoning tasks pose challenges for controlled assessments, particularly in quantifying error propagation and addressing context limitations. Ensuring consistent difficulty across steps is therefore crucial for transparent and reliable evaluation of the reasoning capabilities of large language models, which remains significantly underexplored.

2.3 RULE SYSTEM

Complexity emerges from simple rules. Modern scientific subjects such as mathematics, physics, and chemistry are generally founded on rule-based systems. These rely on axioms Suppes (2012),

162 theorems Russell (2020), and laws Feynman (1963) to guide reasoning and model phenomena.
 163 Arithmetic is the most basic example, providing core operations and relations that support more
 164 advanced structures. In physics, Newton’s laws of motion describe how forces affect bodies, forming
 165 the basis of classical mechanics. Cellular automata illustrate how simple local rules can give rise to
 166 complex global behavior Conway et al. (1970). Rule system is also essential in program verification
 167 and automated reasoning, facilitating formal proofs through Lean theorem prover Moura & Ullrich
 168 (2021), for example. We define the ability of an LLM to faithfully execute rules step by step as
 169 computational reasoning. Based on this definition, we propose TMBench, a rule-based system
 170 grounded in Turing machines. In TMBench, LLM is prompted to output the result of each rule
 171 execution step, which is then evaluated against ground-truth trajectories for assessment. Accurate
 172 intermediate steps are essential to the integrity of the rule-based system.

174 3 TURING MACHINE BENCH

176 3.1 MOTIVATION

178 Reasoning is one of the most important abilities of LLMs and has gain significant attention recently.
 179 Current reasoning benchmarks focus on the final results and depend on domain-specific knowledge.
 180 For instance, the AIME benchmark, used to evaluate mathematical reasoning abilities, relies on
 181 the mastery of mathematical knowledge. Furthermore, there is a risk that models might tailor their
 182 responses to specific benchmark metrics, thereby inflating their performance. From an intuitive
 183 perspective, reasoning is the process of **selecting** appropriate rules and **applying** them faithfully
 184 step by step. We refer to this capability as *computational reasoning*. Computational reasoning is
 185 the ability to systematically select and accurately apply rules, ensuring that each step is transparent,
 186 verifiable, and grounded within the given rule system. In contrast to logical reasoning, which primarily
 187 focuses on the validity of conclusions, computational reasoning emphasizes the faithfulness and
 188 traceability of the entire reasoning process. Therefore, we propose an evaluation framework to
 189 evaluate computational reasoning by Universal Turing Machine Simulation.

190 3.2 TURING MACHINE

192 **Turing Machine (TM)** is a foundational model of computation Turing et al. (1936), defining algo-
 193 rithmic processes and computational limits. Formally, a TM is represented as $\mathcal{T} = (Q, \Gamma, \delta, q_0, H)$,
 194 where Q is the set of states, Γ the tape alphabet, $\delta : Q \times \Gamma \rightarrow Q \times \Gamma \times \{L, R\}$ the transition function,
 195 q_0 the initial state, and H the set of halting states. A TM manipulates symbols on an infinite tape
 196 using a read-write head, guided by a transition function δ , and executes computations step by step
 197 until it reaches a halting state. It serves as a fundamental abstraction for formalizing decidability
 198 and complexity classes. Beyond theoretical significance, TM are valuable for modeling symbolic
 199 reasoning and structured inference. In this work, we leverage Turing Machine principles to construct
 200 a rigorous multi-step benchmark, evaluating LLMs’ ability to perform structured computations which
 201 is vital for reasoning.

202 **Universal Turing Machine (UTM)** is a theory model of computation, capable of simulating the
 203 behavior of any other Turing machine. As a foundation of computability theory, the UTM formalizes
 204 the notion of computational universality and provides the theoretical foundation for the Church-Turing
 205 thesis, which asserts that any effectively computable function can be executed by a Turing Machine.
 206 By encoding both a machine’s description and its input on its tape, a UTM demonstrates that a
 207 single device can emulate any computational process, establishing the basis for general-purpose
 208 computation and modern computing architectures.

209 To facilitate theoretical analysis while preserving computational expressiveness, we adopt the *m*-Tag
 210 System, a simple and recognized model. The *m*-Tag System has been rigorously proven to be Turing-
 211 complete for $m > 1$ Wang et al. (1971); Cocke & Minsky (1964), making it a suitable abstraction for
 212 modeling universal computation within a structured and analyzable framework.

213 **m-Tag System** is a formal computational model introduced in Post (1943) as a simplified yet
 214 computationally equivalent variant of the Universal Turing Machine. It operates on a queue of
 215 symbols, iteratively applying production rules to modify the sequence. A tag system is formally
 described by a triplet (m, A, P) , where:

- m is the deletion number, specifying how many symbols are removed from the head of queue per step.
- A is a finite alphabet of symbols, from which queue elements are drawn.
- P is a set of production rules mapping each $x \in A$ to a corresponding word $P(x)$, which is appended to tail of the queue.

The single-step process is formally defined as follows:

$$\text{Step} : \underbrace{x_1 x_2 \dots x_m}_{\text{Read}} X \longrightarrow \underbrace{x_1 x_2 \dots x_m}_{\text{Delete}} X \underbrace{P(x_1)}_{\text{Write}}, \quad (1)$$

where new symbols $P(x_1)$ are appends to the tail of queue generated based on the head symbol x_1 and production rule. Simultaneously, m symbols are deleted from the head of the queue. This process resembles the next-token prediction mechanism employed by LLMs.

Tag systems are proven Turing completeness when $m > 1$, making them the minimal yet effective computational models, see Section A for complete certification. An example of a 2-tag system simulation is provided in Table 1, illustrating the iterative process of reading, writing, and deletion until a halting condition is reached.

Table 1: Example of 2-tag systems simulation. At each step, the head symbol of the queue is read, and new symbols are appended to the tail based on production rules. The first two symbols are then removed. The system halts when the queue contains fewer than two symbols.

Alphabet	{A, B, C, D, E}	{1, 2, 3, 4, 5}	{@, #, \$, %, &}
Init	[B A E E C]	[5 2 3 2]	[\$ @ @ #]
P-Rules	A : E D A B C B : D C : E E E D D D : B C E : D	1 : 5 5 2 2 : 4 2 5 1 3 3 : 4 3 1 4 : 3 4 5 : 3	@ : % \$ # \$ # : & \$: & % : # # % % & : % # & \$
Steps	0. [B A E E C] (Init) 1. [B A E E C D] 2. [E E C D D] 3. [C D E E E D D] ... 16. [B C D D] 17. [D D B C] 18. [B C D] (Halt)	0. [5 2 3 2] (Init) 1. [5 2 3 2 3] 2. [3 2 3 4 3 1] 3. [3 4 3 1 4 3 1] ... 28. [5 1 3343455243134552343] 29. [3 3 43455243134552343431] 30. [4 3 45524313455234343134]	0. [\$ @ @ #] (Init) 1. [\$ @ @ # &] 2. [@ # & % \$ # \$] 3. [& % \$ # \$ # & \$] ... 28. [# & \$ % # & \$ &] 29. [\$ % # & \$ & &] 30. [# & \$ & & &]

3.3 EVALUATION METRICS

To comprehensively evaluate the multi-step instruction following capability of LLMs in the reasoning process, we define three key metrics: Step Accuracy, Step-Weighted Accuracy, and Pass Rate.

Step Accuracy. It quantifies the proportion of correctly predicted queues at a given step in the reasoning process, providing a fine-grained evaluation of a model’s stepwise performance. It is particularly valuable for multi-step reasoning tasks, where errors can accumulate and propagate, significantly impacting overall accuracy. The accuracy at step i is defined as:

$$\text{ACC}(i) = \frac{N_{\text{correct}}(i)}{N_{\text{total}}(i)}, \quad (2)$$

where $N_{\text{correct}}(i)$ and $N_{\text{total}}(i)$ denote the number of correct predictions and the total number of predictions at step i , respectively.

Step-Weighted Accuracy (SWA) To evaluate model performance across a sequence of reasoning steps with an emphasis on later steps, we define the *Step-Weighted Accuracy* as a weighted average of per-step accuracies:

$$\text{SWA} = \frac{1}{\sum_{i=1}^T w_i} \sum_{i=1}^T w_i \cdot \text{ACC}(i), \quad (3)$$

where w_i is the weight assigned to step i . Setting $w_i = 1$ corresponds to uniform weighting across all steps, while $w_i = i$ linearly increases the emphasis on later steps, thereby prioritizing accurate performance in deeper stages of reasoning. We evaluate performance using Step-Weighted Accuracy (SWA) under both uniform (Uni.) and linear (Lin.) weighting schemes.

Pass Rate. The pass rate quantifies the likelihood of successfully completing a given process without errors before termination or reaching the maximum number of allowed steps. This metric focuses more on the final outcome, whereas the previous metrics emphasize the process.

4 EXPERIMENTS AND RESULTS

In this section, we present a series of experiments conducted on TMBench using a diverse set of models and analytical approaches. The experiments and their results are outlined as follows. First, we introduce the experimental setup including datasets and LLMs. Next, we evaluate a broad range of recent open-source LLMs, spanning from 0.6B to 671B parameters, alongside proprietary models such as Gemini, Grok, and Claude. We then assess the correlation between the TMBench Pass Rate and real-world benchmarks. Finally, we conduct ablation studies including unbounded-step, temperature, alphabet, and difficulty ablations to analyze TMBench’s characteristics, which inevitably fail with increasing steps due to its autoregressive nature, remain robust across different alphabets, and allow continuous and scalable control of task difficulty.

4.1 EXPERIMENTAL SETUP

We provide a brief overview of the experimental setup, including the datasets and the latest large language models evaluated. Detailed information is provided in the released code and data.

Datasets. Based on the dataset design methodology introduced in Section 3, we sample 100 instances of m -tag systems with an alphabet size of 5, where $m = 2$. For each system, the rule lengths range from 1 to 5, and the initial string lengths range from 2 to 9. The maximum simulation length is 30, with 11 cases halting early upon meeting the termination condition. Experiments under other parameter settings can be found in Section 4.4. In addition, we evaluate model performance on several established benchmarks, including AIME2024 MAA (2024), MATH500 Lightman et al. (2023), GPQA Diamond Rein et al. (2024), and MMLU Pro Wang et al. (2024b). Further details can be found in the Appendix.

LLMs. We evaluate a diverse set of state-of-the-art LLMs, including open-source models ranging from 0.6B to 671B parameters, as well as proprietary API-only models, covering various model families and architectures to ensure a comprehensive analysis. Our selection includes the instructed versions of the LLaMA models, specifically the 1B, 8B, and 70B variants, released by the LLaMA team Dubey et al. (2024). From the Qwen family, we evaluate both the instruct and preview versions of the 0.6B, 1.7B, 4B, 8B, 14B, 32B, and 72B models Yang et al. (2024), along with the language model component of the QVQ-72B-Preview multimodal system Wang et al. (2024a); Team (2024), as well as two MoE models: Qwen3-30B-A3B and Qwen3-235B-A22B. For the Gemma family, we include both Gemma3-12B and Gemma3-27B Team et al. (2025). We also incorporate the QwQ-32B-Preview and Sky-T1-32B-Preview model Team (2025). For proprietary models, including Gemini-2.5-ProDeepMind (2025), Grok-3 xAI (2025), Claude-3.7-Sonnet-Thinking Anthropic (2025), Doubao-1.5-Pro Doubao (2025), and GPT-4.1 OpenAI (2024). For efficient and high-throughput inference, we leverage the Transformers library Wolf et al. (2020) alongside vLLM Kwon et al. (2023), which provides optimized execution for large-scale model evaluations. Unless otherwise specified, we employ a greedy decoding strategy and set the maximum generation length to 16384 tokens.

Table 2: Performance of large language models across different scales (0.6B to 605B) and API-only models, evaluated on TMBench. Asterisks (*) denote widely recognized reasoning benchmarks.

Model	SWA (Uni.)	SWA (Lin.)	Pass Rate	AIME*	MATH*	GPQA*
0.6B - 8B						
Qwen3-0.6B	0.9	0.1	0	6.7	65.0	21.7
Llama-3.2-1B	3.3	0.3	0	0	18.4	12.6
Qwen3-1.7B	1.5	1.2	1	26.7	86.4	32.8
Qwen3-4B	7.4	3.5	6	56.7	91.6	49.0
Llama-3.1-8B	4.5	0.4	1	0	46.4	21.2
Qwen3-8B	7.2	4.4	8	60.0	91.4	56.6
12B+						
Gemma-3-12B	7.7	1.6	0	13.3	83.2	41.9
Qwen2.5-14B	10.3	2.6	0	16.7	79.2	43.9
R1-Distill-Qwen-14B	29.5	18.0	10	60.0	87.0	53.5
Qwen3-14B	13.8	9.3	2	76.7	93.0	64.1
27B+						
Gemma-3-27B	26.1	12.2	5	40.0	88.8	46.0
Qwen2.5-32B	18.2	6.0	7	10.0	80.6	48.0
Sky-T1-32B-Preview	20.9	9.1	8	40.0	86.2	51.0
QwQ-32B-Preview	19.8	9.8	4	53.3	90.4	62.1
R1-Distill-Qwen-32B	33.5	22.5	10	72.6	94.3	62.1
Qwen3-32B	10.5	8.3	7	70.0	94.8	63.6
Qwen3-30B-A3B	56.7	46.8	16	80.4	95.9	65.8
70B+						
Llama-3.3-70B	45.2	34.3	12	20.0	75.8	43.9
QVQ-72B-Preview	14.4	6.4	3	33.3	82.8	52.5
Qwen2.5-Math-72B	42.1	26.6	15	20.0	85.2	46.0
Qwen2.5-72B	45.6	29.6	19	13.3	82.8	45.5
Llama-4-Scout	7.7	1.2	1	28.3	84.4	58.7
Llama-4-Maverick	39.1	27.8	19	39.0	88.9	67.1
DeepSeek-V3	87.4	84.9	82	52.0	94.2	65.5
DeepSeek-R1	72.2	63.8	45	79.8	97.3	71.5
Qwen3-235B-A22B	45.2	32.6	22	85.7	93.0	70.0
API only						
Qwen-2.5-Max	27.0	11.8	7	23.3	83.5	58.7
OpenAI-O1-mini	37.0	21.6	11	63.6	90.0	60.0
Gemini-1.5-Pro	40.3	18.4	10	20.0	84.2	56.6
OpenAI-O3-mini	37.1	30.5	7	87.3	98.5	79.7
GPT-4.1	58.7	45.7	26	48.1	91.3	66.3
Doubao-1.5-Pro	76.9	69.1	54	33.3	88.6	65.0
Claude-3.7-Sonnet	85.1	78.1	69	61.3	96.2	78.2
Grok-3-Beta	94.6	92.4	86	83.9	94.8	80.2
Gemini-2.5-Pro	96.6	96.2	94	92.0	98.9	84.0

4.2 EVALUATION OF COMPUTATIONAL REASONING

We evaluate the computation reasoning capabilities of the latest LLMs by analyzing their performance across three metrics: *SWA Uniform*, *SWA Linear*, *Pass Rate*. These metrics offer a comprehensive evaluation of the models’ ability to faithfully execute multi-step reasoning. Detailed results are summarized in Table 2, and the accuracy–step curves are shown in Figure 1. We observe that Gemini-2.5-Pro exhibits robust accuracy in computational reasoning, achieving over 90% accuracy at 30 steps, which suggests an emergent ability to simulate Turing Machine. It is found that models smaller than 4B struggle with even the first step. This observation further supports the phenomenon of **emergence** Wei et al. (2022a). The *QVQ* multimodal model exhibits a notable decline, potentially attributable to its multimodal training process.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

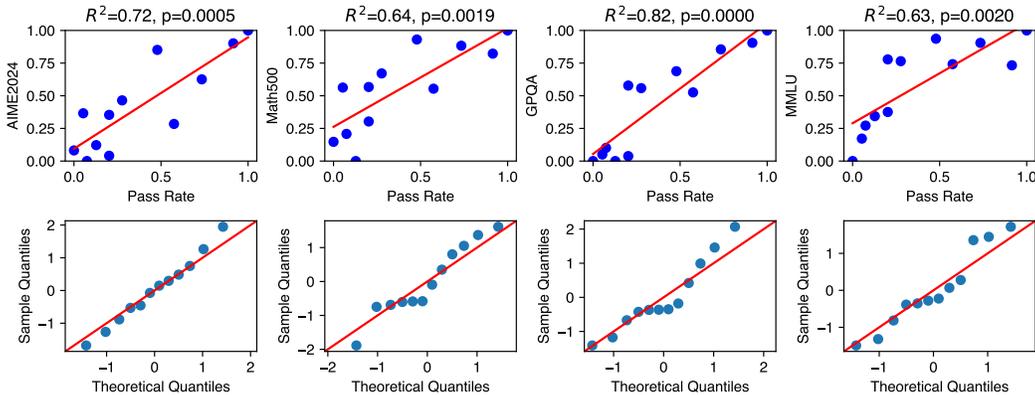


Figure 2: Correlation between TMBench Pass Rate (computational reasoning) and established benchmarks (AIME2024, MATH500, GPQA Diamond and MMLU Pro), where both metrics are min-max normalized, across 12 leading LLMs. Top: Scatter plots with linear fit between TMBench and each benchmark. Bottom: Q-Q plots of regression residuals.

Relationship between computational reasoning and long-context. Figure 1 shows that for models larger than 14B, such as DeepSeek-R1, accuracy decreases linearly as the number of simulated steps grows. This implies a stable per-step error probability, i.e., the difference remains constant with increasing context length. This indicates that for models larger than 14B, as the context increases, their reasoning ability does not decline and the error rate remains stable.

4.3 CORRELATION WITH ESTABLISHED BENCHMARKS

To investigate the relationship between computational reasoning and real-task performance, we analyzed correlations between TMBench Pass Rate and established benchmarks (AIME2024, MATH500, GPQA Diamond, MMLU Pro) across 12 leading LLMs, as shown in Figure 2. The evaluated models include Gemini-2.5-Pro, Grok-3-beta, Claude-3-7-Sonnet, Doubao-1-5-Pro, DeepSeek-R1, GPT-4.1, Llama-4-Maverick, Qwen2.5-72B, Llama-3.3-70B, Qwen2.5-32B, Gemma-3-27B, and Qwen2.5-14B. The results reveal statistically significant correlations ($p < 0.05$), confirming computational reasoning as a predictor of real-world performance. The correlation strength follows the order: GPQA > AIME2024 > MATH500 > MMLU, indicating that tasks requiring deeper reasoning align more closely with TMBench. In contrast, knowledge-heavy benchmarks like MMLU show weaker correlation. The Q-Q plots in the bottom panels suggest residuals are approximately normal, supporting regression assumptions. Furthermore, the average score across AIME2024, MATH500, and GPQA Diamond yields a Pearson correlation coefficient of 0.882 with TMBench (see Figure 5), underscoring the strong link between advanced reasoning and computational ability.

4.4 ABLATION STUDY

Unbounded-Step Ablation. Gemini-2.5-pro demonstrates strong instruction-following capabilities, achieving a 94% pass rate under a 30-step constraint. To further analyze the upper bound of its instruction-following ability, we conducted an unbounded-step evaluation. Using rejection sampling, we selected 10 samples and ran each for up to 1000 steps. Through budget forcing Muennighoff et al. (2025), Gemini was prompted to continue generating until failure with maximum token limit of 1,048,576 tokens. The earliest failure occurred at step 16, and the latest at step 683, as shown in Figure 3a. As an autoregressive model, Gemini inevitably fails with increasing steps due to its statistical nature, underscoring Gemini’s computational limits.

Temperature ablation. To examine the effect of temperature on model performance, we conduct an ablation study with temperature values ranging from 0 to 3.0, while setting the top-p value to 0.8, as shown in Figure 3b. Performance remains stable between 0 and 1, but degrades at higher temperatures.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

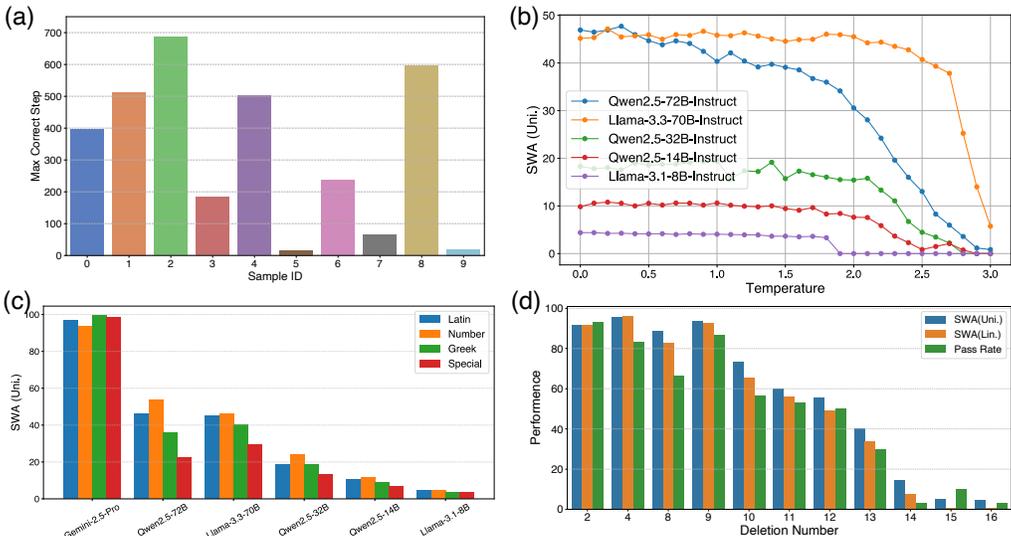


Figure 3: Illustration of ablation results. (a) Maximum correct step achieved by Gemini under unbounded-step execution. (b) Impact of decoding temperature on performance. (c) Effect of alphabet types (Roman, Number, Greek, and Special) on model performance. (d) Task difficulty ablation by varying the number of deletions.

Alphabet ablation. To evaluate the robustness of large language models (LLMs) across diverse alphabet types, we conducted evaluations using four distinct character sets: **Roman letters** (e.g., *a, b, c, d, e*), **Numerals** (e.g., *1, 2, 3, 4, 5*), **Greek letters** (e.g., $\alpha, \beta, \gamma, \delta, \epsilon$), and **Special characters** (e.g., *@, #, \$, %, &*). As illustrated in Figure Figure 3c, while most models exhibit performance fluctuations depending on the character set, Gemini-2.5-Pro consistently achieves high accuracy across all categories, including the *Special characters*, which typically present greater challenges. This robustness indicates that Gemini-2.5-Pro relies on underlying reasoning mechanisms rather than superficial statistical correlations.

Difficulty Ablation. For TMBench, Gemini-2.5-pro achieves a 94% pass rate within 30 steps, with successful simulations extending up to 686 steps, demonstrating its strong computational reasoning capabilities. To further assess this ability, we perform a difficulty ablation study with varying values of *m*, as shown in Figure 3d. The model maintains stable performance when the deletion number *m* is between 2 and 9, but performance declines steadily beyond *m* = 10, approaching zero after *m* = 15. This demonstrates that varying *m* produces a smooth difficulty gradient, underscoring TMBench’s scalability and effectiveness as a benchmarking tool.

5 CONCLUSION

We propose computational reasoning as a fundamental ability of LLMs—the capacity to strictly follow rules and manage internal states for multi-step processes, independent of domain knowledge. To evaluate this, we introduce **Turing Machine Bench (TMBench)**, a benchmark based on m-Tag System Simulation. TMBench is designed with four key attributes: (1) self-contained, knowledge-independent reasoning; (2) interpretable, minimal multi-step structure; (3) adjustable task difficulty; and (4) computational generality grounded in Turing completeness. We evaluate state-of-the-art open-source models (0.6B–671B) and proprietary models such as Gemini and Grok. Results show that TMBench captures computational reasoning effectively, revealing clear multi-step performance curves. Importantly, TMBench aligns well with real-world reasoning benchmarks (e.g., AIME, GPQA), more so than with knowledge-heavy tasks like MMLU. Ablation studies (step limits, temperature, alphabet, difficulty variations, and SFT) further demonstrate TMBench’s robustness. For example, Gemini shows degradation with longer sequences due to its autoregressive nature, yet remains stable with symbol variations—suggesting reliance on true reasoning mechanisms rather than superficial correlations.

REFERENCES

- 486
487
488 Anthropic. Claude 3.7 sonnet thinking, 2025. URL <https://www.anthropic.com/news/claude-3-7-sonnet>. Accessed: 2025-02-25.
489
- 490 Konstantine Arkoudas. Gpt-4 can’t reason. *arXiv preprint arXiv:2308.03762*, 2023.
491
- 492 Zhangir Azerbayev, Bartosz Piotrowski, Hailey Schoelkopf, Edward W Ayers, Dragomir Radev, and
493 Jeremy Avigad. Proofnet: Autoformalizing and formally proving undergraduate-level mathematics.
494 *arXiv preprint arXiv:2302.12433*, 2023.
495
- 496 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain,
497 Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with
498 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 499 BIG bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities of
500 language models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL
501 <https://openreview.net/forum?id=uyTL5Bvosj>.
502
- 503 Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research
504 with large language models. *Nature*, 624(7992):570–578, 2023.
505
- 506 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
507 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,
508 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,
509 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,
510 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios
511 Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino,
512 Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,
513 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa,
514 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob
515 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating
516 large language models trained on code. 2021a.
- 517 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
518 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
519 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021b.
- 520 Cheng-Han Chiang and Hung-yi Lee. Can large language models be an alternative to human
521 evaluations? In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the*
522 *61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,
523 pp. 15607–15631, Toronto, Canada, July 2023. Association for Computational Linguistics. doi:
524 10.18653/v1/2023.acl-long.870. URL [https://aclanthology.org/2023.acl-long.](https://aclanthology.org/2023.acl-long.870)
525 870.
- 526 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
527 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
528 math word problems. *arXiv preprint arXiv:2110.14168*, 2021a.
529
- 530 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
531 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
532 math word problems. *arXiv preprint arXiv:2110.14168*, 2021b.
- 533 John Cocke and Marvin Minsky. Universality of tag systems with $p = 2$. *J. ACM*, 11(1):15–20,
534 January 1964. ISSN 0004-5411. doi: 10.1145/321203.321206. URL [https://doi.org/10.](https://doi.org/10.1145/321203.321206)
535 1145/321203.321206.
536
- 537 John Conway et al. The game of life. *Scientific American*, 223(4):4, 1970.
538
- 539 Google DeepMind. Gemini 2.5 pro, 2025. URL <https://deepmind.google/technologies/gemini/>. Accessed: 2025-03-28.

- 540 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
541 bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of*
542 *the North American chapter of the association for computational linguistics: human language*
543 *technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.
- 544 Doubao. Doubao 1.5 pro, 2025. URL [https://team.doubao.com/en/special/
545 doubao_1_5_pro](https://team.doubao.com/en/special/doubao_1_5_pro). Accessed: 2025-01-22.
- 547 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
548 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
549 *arXiv preprint arXiv:2407.21783*, 2024.
- 550 Richard Phillips Feynman. The feynman lectures on physics. (*No Title*), 1:46, 1963.
- 551
- 552 Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle
553 use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of*
554 *the Association for Computational Linguistics*, 9:346–361, 2021.
- 555 Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Wenfei Zhou, James
556 Coady, David Peng, Yujie Qiao, Luke Benson, et al. Folio: Natural language reasoning with
557 first-order logic. *arXiv preprint arXiv:2209.00840*, 2022.
- 558 Divij Handa, Pavel Dolin, Shrinidhi Kumbhar, Tran Cao Son, and Chitta Baral. Actionreason-
559 ingbench: Reasoning about actions with and without ramification constraints. *arXiv preprint*
560 *arXiv:2406.04046*, 2024.
- 561
- 562 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
563 Steinhardt. Measuring massive multitask language understanding. In *International Conference on*
564 *Learning Representations*, 2020.
- 565 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
566 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*
567 *preprint arXiv:2103.03874*, 2021.
- 568 Yiming Huang, Zhenghao Lin, Xiao Liu, Yeyun Gong, Shuai Lu, Fangyu Lei, Yaobo Liang, Yelong
569 Shen, Chen Lin, Nan Duan, et al. Competition-level problems are effective llm evaluators. *arXiv*
570 *preprint arXiv:2312.02143*, 2023.
- 571
- 572 Zhen Huang, Zengzhi Wang, Shijie Xia, Xuefeng Li, Haoyang Zou, Ruijie Xu, Run-Ze Fan, Lyuman-
573 shan Ye, Ethan Chern, Yixin Ye, et al. Olympicarena: Benchmarking multi-discipline cognitive
574 reasoning for superintelligent ai. *arXiv preprint arXiv:2406.12753*, 2024.
- 575 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
576 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint*
577 *arXiv:2310.06770*, 2023.
- 578
- 579 Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly
580 supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.
- 581 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
582 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
583 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating*
584 *Systems Principles*, 2023.
- 585 Emanuele La Malfa, Christoph Weinhuber, Orazio Torre, Fangru Lin, Samuele Marro, Anthony Cohn,
586 Nigel Shadbolt, and Michael Wooldridge. Code simulation challenges for large language models.
587 *arXiv preprint arXiv:2401.09074*, 2024.
- 588
- 589 Junyu Lai, Jiahe Xu, Yao Yang, Yunpeng Huang, Chun Cao, and Jingwei Xu. Executing arithmetic:
590 Fine-tuning large language models as turing machines. *arXiv preprint arXiv:2410.07896*, 2024.
- 591 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-
592 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative
593 reasoning problems with language models. *Advances in Neural Information Processing Systems*,
35:3843–3857, 2022.

- 594 Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy
595 Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following
596 models. https://github.com/tatsu-lab/alpaca_eval, 2023.
597
- 598 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom
599 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation
600 with alphacode. *Science*, 378(6624):1092–1097, 2022.
- 601 Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yan
602 Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language
603 models. *arXiv preprint arXiv:2211.09110*, 2022.
604
- 605 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
606 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
607 *arXiv:2305.20050*, 2023.
- 608 Zi Lin, Zihan Wang, Yongqi Tong, Yangkun Wang, Yuxin Guo, Yujia Wang, and Jingbo Shang.
609 ToxicChat: Unveiling hidden challenges of toxicity detection in real-world user-AI conversa-
610 tion. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for*
611 *Computational Linguistics: EMNLP 2023*, pp. 4694–4702, Singapore, December 2023. As-
612 sociation for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.311. URL
613 <https://aclanthology.org/2023.findings-emnlp.311>.
- 614 Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A
615 challenge dataset for machine reading comprehension with logical reasoning. *arXiv preprint*
616 *arXiv:2007.08124*, 2020.
617
- 618 Nicholas Lourie, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Unicorn on rainbow: A
619 universal commonsense reasoning model on a new multitask benchmark. In *Proceedings of the*
620 *AAAI Conference on Artificial Intelligence*, volume 35, pp. 13480–13488, 2021.
- 621 Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,
622 Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for
623 science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521,
624 2022.
625
- 626 Man Luo, Shrinidhi Kumbhar, Mihir Parmar, Neeraj Varshney, Pratyay Banerjee, Somak Aditya,
627 Chitta Baral, et al. Towards logigluue: A brief survey and a benchmark for analyzing logical
628 reasoning capabilities of language models. *arXiv preprint arXiv:2310.00836*, 2023.
- 629 MAA. American invitational mathematics examination - aime. In
630 *American Invitational Mathematics Examination - AIME 2024*, Febru-
631 ary 2024. URL [https://maa.org/math-competitions/](https://maa.org/math-competitions/american-invitational-mathematics-examination-aime)
632 [american-invitational-mathematics-examination-aime](https://maa.org/math-competitions/american-invitational-mathematics-examination-aime).
633
- 634 Leonardo de Moura and Sebastian Ullrich. The lean 4 theorem prover and programming language. In
635 *Automated Deduction—CADE 28: 28th International Conference on Automated Deduction, Virtual*
636 *Event, July 12–15, 2021, Proceedings 28*, pp. 625–635. Springer, 2021.
- 637 Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke
638 Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time
639 scaling. *arXiv preprint arXiv:2501.19393*, 2025.
640
- 641 Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David
642 Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work:
643 Scratchpads for intermediate computation with language models. *arXiv preprint arXiv:2112.00114*,
644 2021.
- 645 OpenAI. Openai o1 system card, Sep. 2024. URL [https://](https://assets.ctfassets.net/kftzwdyauwt9/67qJD51Aur3eIc96i0feOP/71551c3d223cd97e591aa89567306912/o1_system_card.pdf)
646 [assets.ctfassets.net/kftzwdyauwt9/67qJD51Aur3eIc96i0feOP/](https://assets.ctfassets.net/kftzwdyauwt9/67qJD51Aur3eIc96i0feOP/71551c3d223cd97e591aa89567306912/o1_system_card.pdf)
647 [71551c3d223cd97e591aa89567306912/o1_system_card.pdf](https://assets.ctfassets.net/kftzwdyauwt9/67qJD51Aur3eIc96i0feOP/71551c3d223cd97e591aa89567306912/o1_system_card.pdf). (Accessed on
2024/09/23).

- 648 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong
649 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton,
650 Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and
651 Ryan Lowe. Training language models to follow instructions with human feedback, 2022.
- 652
653 Mihir Parmar, Nisarg Patel, Neeraj Varshney, Mutsumi Nakamura, Man Luo, Santosh Mashetty,
654 Arindam Mitra, and Chitta Baral. Towards systematic evaluation of logical reasoning ability of
655 large language models. *arXiv preprint arXiv:2404.15522*, 2024.
- 656 Nisarg Patel, Mohith Kulkarni, Mihir Parmar, Aashna Budhiraja, Mutsumi Nakamura, Neeraj
657 Varshney, and Chitta Baral. Multi-logieval: Towards evaluating multi-step logical reasoning
658 ability of large language models. *arXiv preprint arXiv:2406.17169*, 2024.
- 659
660 Emil L Post. Formal reductions of the general combinatorial decision problem. *American journal of*
661 *mathematics*, 65(2):197–215, 1943.
- 662 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani,
663 Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In
664 *First Conference on Language Modeling*, 2024.
- 665
666 Bertrand Russell. *Principles of mathematics*. Routledge, 2020.
- 667
668 Patrick Suppes. *Axiomatic set theory*. Courier Corporation, 2012.
- 669
670 Oyvind Tafjord, Bhavana Dalvi Mishra, and Peter Clark. Proofwriter: Generating implications,
671 proofs, and abductive statements over natural language. *arXiv preprint arXiv:2012.13048*, 2020.
- 672 Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej,
673 Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical
674 report. *arXiv preprint arXiv:2503.19786*, 2025.
- 675
676 NovaSky Team. Sky-t1: Train your own o1 preview model within \$450. 2025.
- 677 Qwen Team. Qvq: To see the world with wisdom, December 2024. URL [https://qwenlm.](https://qwenlm.github.io/blog/qvq-72b-preview/)
678 [github.io/blog/qvq-72b-preview/](https://qwenlm.github.io/blog/qvq-72b-preview/).
- 679
680 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
681 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
682 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 683
684 Alan Mathison Turing et al. On computable numbers, with an application to the entscheidungsproblem.
685 *J. of Math*, 58(345-363):5, 1936.
- 686
687 Hao Wang, John Cocke, Marvin Minsky, and Stephen A. Cook. Tag systems and lag systems. *Journal*
688 *of Symbolic Logic*, 36(2):344–344, 1971. doi: 10.2307/2270314.
- 689
690 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,
691 Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the
692 world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024a.
- 693
694 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
695 ury, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.
696 *arXiv preprint arXiv:2203.11171*, 2022.
- 697
698 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming
699 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging multi-
700 task language understanding benchmark. In *The Thirty-eight Conference on Neural Information*
701 *Processing Systems Datasets and Benchmarks Track*, 2024b.
- 702
703 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
704 Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models.
705 *arXiv preprint arXiv:2206.07682*, 2022a.

- 702 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
703 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
704 *neural information processing systems*, 35:24824–24837, 2022b.
- 705
706 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
707 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art
708 natural language processing. In *Proceedings of the 2020 conference on empirical methods in*
709 *natural language processing: system demonstrations*, pp. 38–45, 2020.
- 710 Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal
711 llm. *arXiv preprint arXiv:2309.05519*, 2023.
- 712
713 xAI. Grok-3, 2025. URL <https://x.ai/grok>. Accessed: 2025-02-19.
- 714 Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-
715 answering to explaining temporal actions. In *Proceedings of the IEEE/CVF conference on computer*
716 *vision and pattern recognition*, pp. 9777–9786, 2021.
- 717
718 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,
719 Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint*
720 *arXiv:2412.15115*, 2024.
- 721 Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual
722 commonsense reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and*
723 *pattern recognition*, pp. 6720–6731, 2019a.
- 724
725 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine
726 really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for*
727 *Computational Linguistics*, pp. 4791–4800, 2019b.
- 728
729 Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. Minif2f: a cross-system benchmark for
730 formal olympiad-level mathematics. *arXiv preprint arXiv:2109.00110*, 2021.
- 731
732 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
733 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
734 Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on*
735 *Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=uccHPGDlao>.
- 736
737 Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu
738 Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation models.
739 *arXiv preprint arXiv:2304.06364*, 2023.
- 740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756 A THEORY.

757
758 In this section, we demonstrate that the m -tag ($m_i 1$) system constitutes a universal Turing machine,
759 following the construction presented in Cocke & Minsky (1964), with minor corrections.
760

761 A.1 REPRESENTATION OF A TURING MACHINE

762
763 We begin by introducing the definition and step operation for the Turing Machine. Consider a
764 Turing Machine with a binary alphabet $\Gamma = \{0, 1\}$, where each state Q_i is defined by a quadruple
765 $(S_i, D_i, Q_{i0}, Q_{i1})$, corresponding to the transitions for the symbols 0 and 1. A single step of the
766 Turing Machine, when in state Q_i , is performed as follows:

- 767 1. Write the symbol S_i onto the current tape cell.
- 768 2. Move the tape head one cell in direction D_i , where $D_i \in \{L, R\}$.
- 769 3. Read the symbol S' from the cell currently.
- 770 4. Transition to the next state based on the symbol read S' :
 - 771 • If $S' = 0$, the next state is Q_{i0} .
 - 772 • If $S' = 1$, the next state is Q_{i1} .

773
774
775 In this specific transition rule, the symbol written S_i and the tape head direction D_i are determined
776 solely by the current state Q_i . The subsequent state transition depends on the symbol S' read from
777 the tape after the move: the machine transitions to state Q_{i0} if $S' = 0$, and to state Q_{i1} if $S' = 1$.

778 An instantaneous description of the machine with state Q_i at any given step is represented as:

$$779 \dots a_3 a_2 a_1 a_0 \alpha b_1 b_2 b_3 \dots$$

780 where Q_i is the state after reading α . The complete instantaneous description can be described by a
781 *triplet*:

$$782 (Q, M, N) = (Q_i, \sum_{i=0}^{\infty} a_i 2^i, \sum_{i=0}^{\infty} b_i 2^i)$$

783
784 where M is the string of symbols to the left of α , N is the string of symbols to the right of α . The
785 symbol α is the current symbol being read, which will soon be replaced by S_i depending on the
786 transition rule.

787 For a rightward movement, the next state Q' depends on the value of S_i . If $S_i = 0$, then $Q' = Q_{i0}$,
788 and if $S_i = 1$, then $Q' = Q_{i1}$.

789 The updating of the strings follows the rules:

$$790 M \leftarrow 2M + S_i$$

$$791 N \leftarrow \left\lfloor \frac{N}{2} \right\rfloor$$

792
793 In the case of a rightward movement, only the values of M and N are swapped, due to the symmetry
794 of the process. For the sake of simplicity, we assume here that the focus is on the rightward move,
795 without loss of generality.

800 A.2 EQUIVALENT CONSTRUCTION WITH A TAG SYSTEM

801
802 Given the instantaneous description of the Turing Machine as a triplet (Q_i, M, N) , we can construct
803 a corresponding string using the following form:

$$804 A_i x_i (\alpha_i x_i)^M B_i x_i (\beta_i x_i)^N$$

805
806 For simplicity and clarity, we will omit the subscript i in the notation and rewrite it as:

$$807 Ax(\alpha x)^M Bx(\beta x)^N$$

810 STEP 0 (INITIAL)

$$811 \quad Ax(\alpha x)^M Bx(\beta x)^N \quad (4)$$

814 STEP 1

815 Rules:

$$816 \quad A \rightarrow \begin{cases} Cx, & \text{if } S = 0 \\ Cxcx, & \text{if } S = 1 \end{cases}, \quad \alpha \rightarrow cxcx$$

819 Updated:

$$820 \quad Bx(\beta x)^N Cx(cx)^{M'} \quad (5)$$

822 where

$$823 \quad M' = \begin{cases} 2M, & \text{if } S = 0 \\ 2M + 1, & \text{if } S = 1 \end{cases}$$

826 STEP 2

828 Rules:

$$829 \quad B \rightarrow S, \quad \beta \rightarrow s$$

830 Updated:

$$831 \quad Cx(cx)^{M'} S(s)^N \quad (6)$$

834 STEP 3

835 Rules:

$$836 \quad C \rightarrow D_1 D_0, \quad c \rightarrow d_1 d_0$$

838 Updated:

$$839 \quad S(s)^N D_1 D_0 (d_1 d_0)^{M'} \quad (7)$$

841 STEP 4

843 Rules:

$$844 \quad S \rightarrow T_1 T_0, \quad s \rightarrow t_1 t_0$$

845 Based on the parity of N , the updated string takes different forms:

847 If N is odd:

$$848 \quad D_1 D_0 (d_1 d_0)^{M'} T_1 T_0 (t_1 t_0)^{\frac{N-1}{2}},$$

849 If N is even:

$$850 \quad D_0 (d_1 d_0)^{M'} T_1 T_0 (t_1 t_0)^{\frac{N}{2}}.$$

852 STEP 5

853 If N is odd:

$$854 \quad D_1 \rightarrow A_1 x_1, \quad d_1 \rightarrow a_1 x_1$$

855 Then, the resulting string is:

$$856 \quad T_1 T_0 (t_1 t_0)^{\frac{N-1}{2}} A_1 x_1 (\alpha_1 x_1)^{M'} \quad (8)$$

859 If N is even:

$$860 \quad D_0 \rightarrow x_0 A_0 x_0, \quad d_0 \rightarrow a_0 x_0$$

861 Then, the resulting string is:

$$862 \quad T_0 (t_1 t_0)^{\frac{N}{2}} x_0 A_0 x_0 (\alpha_0 x_0)^{M'} \quad (9)$$

863

864 STEP 6 (FINAL)

865 If N is odd:

$$866 T_1 \rightarrow B_1x_1, \quad t_1 \rightarrow \beta_1x_1$$

867 Then,

$$868 A_1x_1(\alpha_1x_1)^{M'} B_1x_1(\beta_1x_1)^{\frac{N-1}{2}} \quad (10)$$

869 If N is even:

$$870 T_0 \rightarrow B_0x_0, \quad t_0 \rightarrow \beta_0x_0$$

871 Then,

$$872 A_0x_0(\alpha_0x_0)^{M'} B_0x_0(\beta_0x_0)^{\frac{N}{2}} \quad (11)$$

873 This completes the proof that a single computation step of a Turing machine is successfully simulated
874 by a sequence of transitions in the 2-tag system.

875 B EXPERIMENT DETAILS

876 B.1 DATASET DETAILS

877 **AIME2024** MAA (2024) is a collection of problems from the 2024 American Invitational Mathematics Examination. It contains 30 high-school level mathematics problems requiring creative problem solving and deep mathematical insight. These problems are typically used to select top students for the USA Mathematical Olympiad. The dataset emphasizes algebraic manipulation, mathematical reasoning, and non-routine problem solving.

878 **Math500** Lightman et al. (2023) is a representative benchmark consisting of 500 mathematics problems, sampled at random from the MATHHenrycks et al. (2021) evaluation dataset by OpenAI. It includes problems of difficulty across topics such as algebra, geometry, and number theory.

879 **GPQA Diamond** Rein et al. (2024) is a high-difficulty subset of Graduate-Level Google-Proof Q&A Benchmark (GPQA), consisting of 198 multiple-choice questions in biology, physics, and chemistry. It includes only questions where both experts answer correctly and the majority of non-experts answer incorrectly to ensure highest quality.

880 **MMLU Pro** Wang et al. (2024b) is a professional-level extension of the original Massive Multitask Language Understanding (MMLU) benchmark Henrycks et al. (2020), comprising expert-level questions across 57 domains, including mathematics, history, computer science, law, medicine, engineering, the natural sciences, and more. It enhances the original benchmark by introducing more challenging, reasoning-intensive questions, expanding answer choices from four to ten, and removing trivial or noisy items.

881 **TMBench** summarizes the distributions of Step Length, Rule Length, and Initial Length as shown in the following tables.

901 Table 3: Distribution of Step Length

Step Range	1–6	7–12	13–18	19–24	25–30
Avg.	7.61	13.60	19.41	25.11	34.35

902 Table 4: Distribution of Rule Length

Range	1.40–2.08	2.08–2.76	2.76–3.44	3.44–4.12	4.12–4.80
Percentage	0.08	0.25	0.39	0.23	0.05

903 B.2 IMPLEMENTATION DETAILS

904 **Environment.** Our method is implemented with Python 3.10.16, CUDA 12.4, PyTorch 2.6.0 and vLLM 0.8.5. The required libraries are specified in the `requirements.txt` file provided in the repository. The experiments are performed on a machine with 96 vCPUs (2.90 GHz) from Intel Xeon processors, eight NVIDIA A100-80GB GPUs, and 1024 GB of RAM.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Table 5: Distribution of Initial Length

Range	2.00–3.40	3.40–4.80	4.80–6.20	6.20–7.60	7.60–9.00
Percentage	0.18	0.32	0.28	0.14	0.08

Inference Configuration. For API-only models, we use their official APIs and set a maximum token length limit accordingly. For open-source models that exceed the GPU memory limits of a single machine, such as LLaMA 4 Maverick, DeepSeek-V3, and DeepSeek-R1, we rely on APIs provided by Alibaba Cloud and DeepSeek. For all other models, including Qwen, LLaMA, and Gemma, we use vLLM as the inference backend, which is highly efficient. We apply a sampling strategy based on greedy decoding with a maximum token limit of 16,384.

Prompt. We adopt the following 1-shot reasoning prompt template for evaluation, where the shot example is used to guide the output format extraction.

```

972 <|im_start|>user
973 Simulate a m-tag system. Your task is to simulate each
974 transition step-by-step and provide the queue's state at each
975 step. Follow the rules and examples closely, and stop upon
976 reaching the halt condition or 30 steps. Do not generate
977 additional examples or new problems. No code.
978 ## Rules for Simulation:
979 1. In each transition, the machine performs the following
980 steps:
981   - If the queue length is less than m, halt
982   - Read the head symbol of queue
983   - Append symbols to the tail based on the head symbol and the
984     corresponding transition rule
985   - Delete m symbols from the head of the queue
986
987 2. The machine halt if:
988   - The queue's length is less than m.
989
990 ## Example:
991 m: 2
992 Alphabet: A, B, C
993 Init: [B C A]
994 Transition rules:
995 A : C A C
996 B : A
997 C : B
998 Simulation steps:
999 ### step 0:
1000 - Action: Init
1001 - Queue State: [B C A]
1002
1003 ### step 1:
1004 - Head Symbol: B
1005 - Action: Append A to the end of the queue. Remove B C from
1006   the head.
1007 - Queue State: [A A]
1008   ...
1009 ### step 4:
1010 - Head Symbol: C
1011 - Action: Append B to the end of the queue. Remove C B from
1012   the head.
1013 - Queue State: [B] <halt>
1014
1015 ---
1016 ## The Only Problem to Solve:
1017 m: {m}
1018 Alphabet: {ALPHABET}
1019 Init: {INIT}
1020 Transition Rules:
1021 {RULES}
1022 Simulation steps:
1023 <|im_end|>
1024 <|im_start|>assistant
1025

```

1026 **Case.** The response generated by Gemini-2.5-Pro, showcasing its impressive capabilities in instruc-
 1027 tion following and reasoning.
 1028

```

1029 <|im_start|>assistant
1030 Here is the step-by-step simulation:
1031 m: 2
1032 Alphabet: A, B, C, D, E
1033 Init: [B D D]
1034 Transition Rules:
1035 A : C
1036 B : E C E C
1037 C : B B B A
1038 D : D B B
1039 E : A E E E
1039 Simulation steps:
1040 ### step 0:
1041 - Action: Init
1042 - Queue State: [B D D]
1043 ### step 1:
1044 - Head Symbol: B
1045 - Action: Append E C E C to the end of the queue. Remove B D
1046 from the head.
1047 - Queue State: [D E C E C]
1048 ### step 2:
1049 - Head Symbol: D
1050 - Action: Append D B B to the end of the queue. Remove D E
1051 from the head.
1052 - Queue State: [C E C D B B]
1053 ...
1053 ### step 18:
1054 - Head Symbol: E
1055 - Action: Append A E E E to the end of the queue. Remove E C
1056 from the head.
1057 - Queue State: [E C A E E E A E E E A E E E A E E E A E E E A
1058 E E E A E E E A E E E A E E E]
1059 ### step 19:
1060 - Head Symbol: E
1061 - Action: Append A E E E to the end of the queue. Remove E C
1062 from the head.
1063 - Queue State: [A E E E A E E E A E E E A E E E A E E E A E E
1064 E A E E E A E E E A E E E - - -] (Answer)
1065 - Queue State: [A E E E A E E E A E E E A E E E A E E E A E E
1066 E A E E E A E E E A E E E A E E E] (Ground Truth)
1067 ...

```

1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

B.3 MORE RESULTS

Arithmetic tasks. To comprehensively evaluate mathematical reasoning, we consider a suite of challenging benchmarks, including AIME, GSM8K, MATH, GPQA, and MMLU, with results summarized in Table 6.

Model	TMBench	GSM8K	AIME	MATH	GPQA	MMLU
Qwen3-14B	30.6	96.3	76.7	93.0	64.1	67.5
Gemma3-27B-IT	25.8	94.7	40.0	88.8	46.0	66.9
Qwen3-30B-A3B	54.3	96.2	80.4	95.9	65.8	71.0
Qwen3-32B	46.7	96.5	70.0	94.8	63.6	72.7
Qwen3-235B-A22B	46.6	96.7	85.7	93.0	70.0	76.2

Table 6: Experimental results on arithmetic and reasoning benchmarks.

Ablation on Alphabet Size. We observe that when other parameters are held constant, models perform worse with a small alphabet size. With a small alphabet, such as size = 2, the same symbols frequently appear in production rules. Different rules tend to look alike creating similar patterns that large language models struggle to differentiate, leading to more errors. As the alphabet size increases, each symbol appears less frequently, making the rules more distinct and diverse. This uniqueness allows the model to more easily manage the internal states to give the correct results (see Table 7).

Alphabet Size	Qwen3-14B	Gemma3-27B-IT	Qwen3-30B-A3B	Qwen3-32B	Qwen3-235B-A22B	Avg
2	13.2	8.7	13.9	20.2	15.4	14.3
5	30.6	25.8	54.3	46.7	46.6	40.8
8	51.1	31.8	75.5	56.6	65.5	56.1
12	64.3	44.5	80.3	77.0	70.9	67.4
18	67.8	52.1	90.4	86.7	84.8	76.4
26	72.7	55.8	90.7	90.1	83.6	78.6

Table 7: Performance comparison across different alphabet sizes.

Ablation on Rule Length. The results in the Table 8 indicate that increasing the rule length leads to higher task complexity and a corresponding decline in model performance.

Rule Length	Qwen3-14B	Gemma3-27B-IT	Qwen3-30B-A3B	Qwen3-32B	Qwen3-235B-A22B	Avg
1-3	41.3	27.9	63.6	64.0	63.3	52.0
1-4	38.6	24.4	58.8	51.4	53.4	45.3
1-5	30.6	25.8	54.3	46.7	46.6	40.8
2-6	28.7	24.9	51.5	44.0	46.9	39.2
2-7	22.2	19.6	49.0	39.0	41.2	34.2

Table 8: Performance comparison across different Rule Length.

Ablation on Initial Length. We conducted ablations across different initial lengths and found that model performance shows no significant correlation with initial length, suggesting that initial length has minimal impact on the overall benchmark difficulty, as show in Table 9

Impact of SFT on TMBench. To investigate whether TMBench primarily rewards surface-pattern learning, we performed supervised fine-tuning on synthetic traces with different alphabet sizes ($A = 2$ and $A = 5$), as well as on TEST configurations. While SFT consistently improved performance across all models, the gains remained limited compared to frontier systems such as Gemini-2.5-Pro, Claude-3, and Grok-2. Moreover, training on TEST traces did not yield substantially higher performance than training on Non-TEST traces. These results indicate that although SFT helps models capture certain superficial patterns, it is insufficient to bridge the performance gap, reinforcing

Initial Length	Qwen3-14B	Gemma3-27B-IT	Qwen3-30B-A3B	Qwen3-32B	Qwen3-235B-A22B	Avg
2-5	28.4	23.1	52.7	44.4	46.5	39.0
2-7	32.3	18.3	53.6	43.6	49.3	39.4
2-9	30.6	25.8	54.3	46.7	46.6	40.8
3-11	40.6	29.0	57.4	48.1	44.1	40.2
5-13	42.4	29.7	52.7	46.8	49.3	40.3

Table 9: Performance comparison across different Initial Length.

our claim that TMBench demands genuine reasoning and long-horizon state tracking rather than simple template memorization.

Table 10 reports the quantitative outcomes. Although SFT improves model performance, the magnitude of improvement is significantly smaller than in other benchmarks, where TEST-based SFT typically delivers substantial gains. The relatively limited benefits on TMBench thus underscore its difficulty and robustness against shortcut learning.

Model	Llama-3.1-8B	Qwen3-8B	Gemma3-12B	Qwen3-14B
Baseline	4.4	22.8	7.7	30.6
SFT ($A = 2$)	5.5	26.4	11.7	32.5
SFT ($A = 5$)	9.4	24.7	28.0	37.8
SFT ($A = 5$, TEST)	9.2	28.5	28.2	39.7

Table 10: Performance of different models on TMBench before and after SFT. All models were fine-tuned with a learning rate of 1.0×10^{-4} for 3 epochs.

Prompt ablation. Since large language models (LLMs) are known to be sensitive to prompt design, we further investigated the effect of prompt variation on TMBench. In our setup, we adopted a 1-shot chain-of-thought (CoT) template to ensure that model outputs follow a step-by-step format suitable for automatic extraction using regular expressions, thereby enabling fine-grained evaluation at the step level. To assess the robustness of model rankings, we performed ablation studies by varying the number of demonstrations in the prompt (1-shot, 2-shot, 4-shot, and 8-shot). As shown in Table 11, the relative rankings of models remained consistent across different prompt settings, suggesting that TMBench performance is not overly dependent on a particular prompt configuration.

Setting	Qwen3-14B	Gemma3-27B-IT	Qwen3-30B-A3B	Qwen3-32B	Qwen3-235B-A22B
1-shot	30.6	25.8	54.3	46.7	46.6
2-shot	35.3	25.6	53.2	44.8	48.5
4-shot	36.0	26.1	54.7	45.2	43.6
8-shot	35.3	25.9	56.1	47.9	46.0

Table 11: Prompt ablation study on TMBench. Performance of different models under varying numbers of in-context demonstrations.

Token Distribution. We analyze the token distribution across 36 models on TMBench, as illustrated in Figure 4. A larger number of generated tokens does not necessarily correlate with better performance. While models that engage in more reasoning often produce longer outputs, this behavior does not always translate into higher accuracy. In some cases, excessive token counts are primarily due to repetitive output, reflecting the well-known repetition problem. Interestingly, state-of-the-art models such as Gemini-2.5-Pro, Grok, and Claude exhibit similar token distributions, suggesting that high-performing systems strike a balance between reasoning depth and output efficiency.

To provide a more comprehensive view, we further examine average token counts, efficiency (defined as the ratio of total tokens to pass rate), and the distribution of token lengths across different ranges. As shown in Table 12, strong models such as Gemini-2.5-Pro and Grok maintain moderate token lengths and high efficiency, whereas weaker models (e.g., Qwen3-30B-A3B, QwQ-32B-Preview, and R1-Distill-Qwen-14B) tend to produce extremely long outputs with low efficiency. These findings

indicate that TMBench not only evaluates reasoning ability but also reflects a model’s capacity to generate concise and effective outputs.

Model	Pass Rate	Avg Token	Efficiency	< 1000	1000–5000	5000–10000	> 10000
Gemini-2.5-Pro	94	1993.6	2120.9	0.09	0.91	0	0
Grok-3-Beta	86	2182.8	2538.1	0.08	0.92	0	0
Claude-3.7-Sonnet	69	1640.3	2377.2	0.15	0.85	0	0
Doubao-1.5-Pro	54	1932.0	3577.8	0.10	0.90	0	0
Deepseek-R1	45	1833.4	4074.2	0.10	0.90	0	0
Qwen3-30B-A3B	16	8500.8	53129.7	0.01	0.11	0.57	0.31
QwQ-32B-Preview	4	5636.8	140920.2	0.02	0.58	0.28	0.12
R1-Distill-Qwen-14B	10	10235.1	102350.8	0	0.08	0.44	0.48
Qwen3-4B	6	13695.6	228259.8	0	0.11	0.15	0.74

Table 12: Token distribution analysis on TMBench. We report pass rate, average token count, efficiency (total tokens divided by pass rate), and the proportion of outputs within different token length ranges.

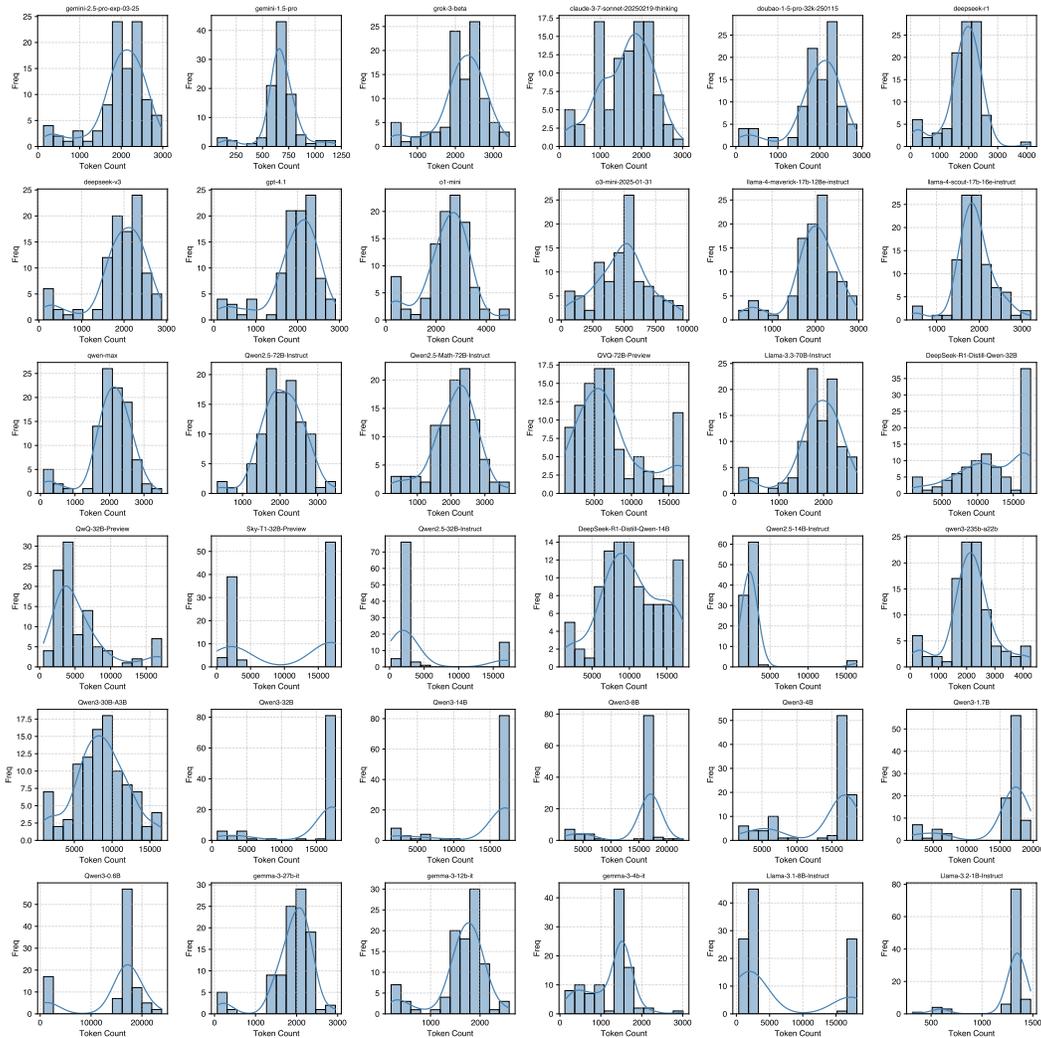


Figure 4: Illustration of token distribution.

Benchmark Correlation. We compute the average score across AIME2024, MATH500, and GPQA Diamond, obtaining a Pearson correlation coefficient of 0.882 with the TMBench pass rate, as

shown in Figure 5. This strong correlation suggests that advanced reasoning performance is closely associated with computational reasoning ability.

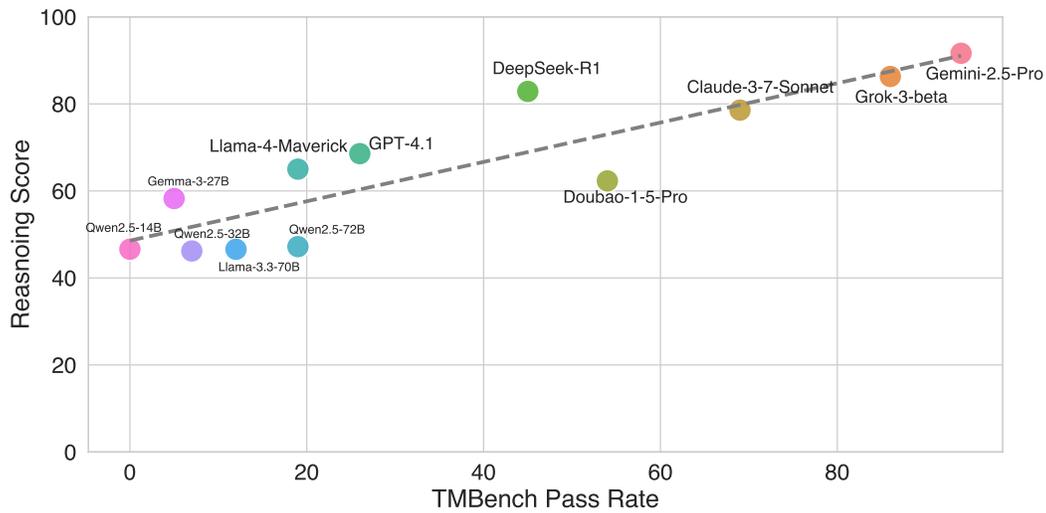


Figure 5: Correlation between TMBench Pass Rate (computational reasoning ability) and Reasoning Score (averaged across AIME2024, MATH500, and GPQA) among LLMs, with a Pearson correlation coefficient of 0.882 with $p=1.49e-04$, demonstrating the connection between TMBench as an abstract rule-based system simulation and real-world reasoning problem.

C LLM IN PAPER WRITING

We use LLMs only to aid and polish writing.