

Do LLMs Really Need 10+ Thoughts for “Find the Time 1000 Days Later”? Towards Structural Understanding of LLM Overthinking

Xinliang Frederick Zhang^{1,2,+}, Anhad Mohananey^{2,*}, Alexandra Chronopoulou^{2,*},
Pinelopi Papalampidi², Somit Gupta², Tsendsuren Munkhdalai²,
Lu Wang¹, and Shyam Upadhyay²

¹University of Michigan, ²Google DeepMind

xlfzhang@{umich.edu, google.com}

**Equal contributions*

Abstract

Models employing long chain-of-thought (CoT) reasoning have shown superior performance on complex reasoning tasks. Yet, this capability introduces a critical and often overlooked inefficiency—**overthinking**—models often engage in unnecessarily extensive reasoning even for *simple* queries, incurring significant computations without accuracy improvements. While prior work has explored solutions to mitigate overthinking, a fundamental gap remains in our understanding of its underlying causes. Most existing analyses are limited to superficial, profiling-based observations, failing to delve into *LLMs’ inner workings*. This study introduces a systematic, fine-grained analyzer of LLMs’ thought process to bridge the gap, TRACE. We first benchmark the overthinking issue, confirming that long-thinking models are five to twenty times slower on simple tasks with no substantial gains. We then use TRACE to first decompose the thought process into minimally complete *sub-thoughts*. Next, by inferring discourse relationships among sub-thoughts, we construct granular thought progression graphs and subsequently identify common thinking patterns for topically similar queries. Our analysis reveals two major patterns for open-weight thinking models—*Explorer* and *Late Landing*. This finding provides evidence that *over-verification* and *over-exploration* are the primary drivers of overthinking in LLMs. Grounded in thought structures, we propose a utility-based definition of overthinking, which moves beyond length-based metrics. This revised definition offers a more insightful understanding of LLMs’ thought progression, as well as practical guidelines for principled overthinking management.

1 Introduction

There is a recent interest in thinking models that generate long-chain-of-thought (CoT) responses

⁺Work done during XFZ’s internship at GDM.

without manual prompts (Xu et al., 2025a; Li et al., 2025). However, current state-of-the-art thinking models (Comanici et al., 2025; Jaech et al., 2024; Zheng et al., 2025; DeepSeek-AI et al., 2025) often engage in excessive thinking even for simple inquiries, leading to computational inefficiency and unnecessary reasoning (Sui et al., 2025). While various methods have been proposed to address this, such as using RL training with length-based penalties (Arora and Zanette, 2025; Aggarwal and Welleck, 2025), compressing reasoning (Deng et al., 2024), or specifying response length in prompts (Han et al., 2025; Xu et al., 2025b), a fundamental knowledge gap remains in our understanding of the underlying causes behind LLM overthinking. Worse still, most of existing analyses are limited to generic profiling and only superficially examine how thinking length impacts model performance (Chen et al., 2024; Yang et al., 2025; Fang et al., 2025). In fact, a deeper analysis of LLMs’ internal mechanisms—specifically, their thought structures and the minimum thinking effort sufficient for effective performance—is conspicuously absent.

To this end, we first systematically benchmark overthinking of LLMs, via head-to-head comparison between *non-thinking* and *thinking* modes.¹ Our results show that for *simple queries* like “7+2” or “date of 1000 days after today”, thinking models naively spend *five to twenty times* longer at inference than non-thinking counterparts, with little or even no improvement in performance. This finding is inline with the literature (Chen et al., 2024) but we uncover the finding at scale, across 14 thinking LLMs (Yang et al., 2025; DeepSeek-AI et al., 2025) and 6 data domains (i.a., Miao et al., 2020; Cobbe et al., 2021; Tan et al., 2023; Wei et al., 2024).

Next, we dig deeper into the LLMs’ thought

¹*Non-thinking* mode provides near-instant responses; *thinking* mode takes longer time to reason before answering.

process to understand the “why” and “how” of overthinking via our proposed analyzer (Figure 2), TRACE, a Thought-process Reconstruction and Automated Clustering Engine. Concretely, we decompose the thought process into smaller, independent, minimally complete sub-thoughts, and then evaluate the accuracy and helpfulness of each sub-thought. Next, we leverage LLM-as-rater (Robertson and Zaragoza, 2009; Comanici et al., 2025) to infer discourse relationships among sub-thoughts. The outputs allow us to construct thought progression graphs for individual queries, and furthermore, identify general thought patterns for *topically similar queries*. Our analysis reveals two primary thought progression patterns: the **Explorer** and **Late Landing**. The *Explorer* pattern distributes the probability of correctness across multiple potential answers. It may find the correct answer early, but it continues to explore alternatives. Conversely, the *Late Landing* pattern follows a convergent reasoning path where the probability of correctness is highly concentrated on the final, distinct answer. This further suggests **over-verification** and **over-exploration** as key drivers of overthinking. Based on the thought structures, we propose a utility-based definition of overthinking—*continuation of the thought process after the marginal return drops below a predefined threshold (ϵ)*. This revised definition moves beyond mere length-based metrics, and transforms overthinking from a vague notion into a quantifiable and actionable metric. Concretely, it not only pinpoints the onset of overthinking within the thought process, but also informs practical heuristics for real-time overthinking management using thought progression-induced characteristics. Overall, TRACE deepens the understanding of LLM overthinking more structurally.

2 Related Work

Thinking Models. Large Language Models (LLMs), built on the transformer architecture (Vaswani et al., 2017), represent a transformative step toward Artificial General Intelligence (AGI; Mumuni and Mumuni, 2025). Recently, the rise of thinking models² (Xu et al., 2025a; Jaech et al., 2024; DeepSeek-AI et al., 2025) has significantly enhanced performance with their **system-2 thinking** competence (Da Silva, 2023; Li et al., 2025). Their remarkable performance on complex tasks, such as challenging mathematical

contests (Hendrycks et al., 2021; Lightman et al., 2024; Luong and Lockhart, 2025), competitive programming tasks (Shi et al., 2024; Jain et al., 2025) and personalization tasks (Zhang et al., 2025a), is largely attributed to their capacity for generating extensive reasoning traces (Li et al., 2025; Chen et al., 2025). This capability is built upon the foundational technique of chain-of-thought (CoT) prompting (Kojima et al., 2022; Wei et al., 2022), which encourages models to articulate intermediate steps before answering. By expending more tokens to explore various problem-solving paths, engaging in self-correction, and performing verification, these models achieve state-of-the-art (SOTA) results in challenging domains (Snell et al., 2024; Muennighoff et al., 2025; Zhang et al., 2025b).

LLM Overthinking. While the “longer is better” paradigm (Jin et al., 2024; Shen et al., 2025), a test-time scaling approach (Snell et al., 2024) that empowers system-2 thinking, has boosted model accuracy on complex tasks, recent research reveals a critical *inefficiency* issue in thinking LLMs: **overthinking** (Chen et al., 2024; Sui et al., 2025). This phenomenon is characterized by models expending *unnecessarily excessive compute* on queries, especially simple ones and those with an evident answer, for marginal improvements or even negative returns. At its worst, overthinking causes performance degradation (Liu et al., 2024), where a model may abandon a correct initial/intermediate answer for a wrong one eventually. Although the benefits of extended reasoning on difficult problems are clear, the trade-offs across the full spectrum of task difficulties remain *underexplored*, particularly for *simple queries*,³ where models’ verbosity makes them more susceptible to overthinking (Chen et al., 2024; Pu et al., 2025). To this end, we are the first to comprehensively benchmark LLM overthinking tendencies on *simple queries*. Our evaluation spans both *horizontally*—across a diverse range of tasks and domains—and *vertically*, by varying the intrinsic difficulty of the problems.

Meanwhile, despite a variety of strategies have been proposed to address the overthinking issue (Arora and Zanette, 2025; Aggarwal and Welleck, 2025; Han et al., 2025; Xu et al., 2025b), a deep exploration of the underlying drivers of overthinking remains largely absent in the literature. Existing analyses are often superficial, focusing on

²They are also commonly referred to as *reasoning models*.

³Though “simple” is cognitively hard to define, we treat queries solvable by bright middle school students as *simple*.

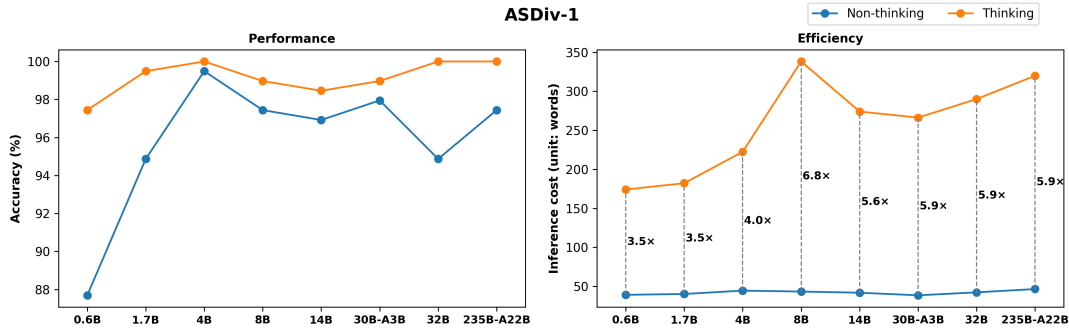


Figure 1: Performance and inference-time efficiency trends of Qwen3 models at different scales on simple math reasoning. We find that additional thinking becomes ineffective once the model scale is above the threshold of 4B. Plots for other tasks such as temporal and logical reasoning, as well as knowledge recall, are in Section C.

coarse, length-based metrics. For instance, some studies measure overthinking in a monolithic way, defining it based on the token distance to the first correct answer and treating subsequent tokens as overthinking (Chen et al., 2024); Others evaluate whether the shortest response among over-sampled generations is the most cost-effective (Pu et al., 2025; Hassid et al., 2025); More recent work, inspired by thinking mode fusion (Yang et al., 2025), quantifies overthinking as the generation of excessive thinking tokens on simple queries without yielding performance improvements over the non-thinking mode (Aggarwal et al., 2025). Regardless of implementation variants, existing approaches remain confined to length-based profiling, without examining the underlying *thought evolution*.

In contrast, our work provides a fine-grained analysis at the sub-thought level. Specifically, we investigate the **internal dynamics of the thought process**, examining how different answers are proposed throughout the thought process and, more importantly, how later sub-thoughts relate to earlier ones. To the best of our knowledge, we are the first to conduct such a structural analysis to understand *why and how* overthinking occurs in LLMs, from its inner workings perspective.

3 Benchmark LLM Overthinking

To investigate thinking LLMs’ inefficiency, we systematically benchmark the phenomenon of overthinking on **simple queries** (Shojaee et al., 2025),⁴ where LLMs expend computational resources on thought tokens that yield marginal or no performance gains. While prior work has identified this issue, its scope has been largely confined to reasoning tasks related to STEM (Aggarwal et al., 2025), especially mathematical reasoning (Chen

⁴The goal of this paper is to provide a focused, controlled study of overthinking. To this end, simple queries allow us to isolate this behavior without confounding factors.

et al., 2024; Hammoud et al., 2025). Our research provides a more comprehensive analysis by evaluating overthinking tendencies both *horizontally* (across various domains) and *vertically* (varying task difficulties). This dual-axis approach provides the first large-scale characterization of the trade-off between reasoning length and model performance, clarifying *when and why thinking is beneficial*.

Defining Overthinking. Inline with the Following (Chen et al., 2024; Sui et al., 2025; Aggarwal et al., 2025), we introduce an *initial* definition of overthinking based on the **generation length**, allowing us to quantify the overthinking extent.

Length-based Overthinking Definition

Overthinking is the generation of additional thought tokens that do not contribute to a performance gain.

Formally, we measure this as the difference in the total generation length between the thinking and non-thinking modes of *the same evaluated model* for the subset of questions that are *already answered correctly* in the non-thinking mode.

3.1 Horizontal Analysis

Our *horizontal* analysis surveys a wide range of domains and tasks—ranging from simple reasoning (e.g., grade-level math problems) to fact-retrieval tasks with minimal reasoning requirements (e.g., machine reading comprehension). The goal is to *establish a baseline by identifying settings where extended reasoning is unlikely to offer benefits*.

Experimental Setup. To ensure our findings are robust, we benchmark a suite of 14 models trained via reinforcement learning (RL) or supervised fine-tuning via distillation. Specifically, we sample responses from a wide range of models, including

Category	Knowledge Recall			Reasoning		
Domain	Short-ctx MRC	Long-ctx Fact Retrieval	Factuality QA	Math	Temporal	Logical
Dataset	SQuAD 2.0	NIAH	SimpleQA	ASDIV	Date Arithmetic	Zebra Logic
Size	420 (210)	352 (132)	400	195	400	240

Table 1: Datasets used for benchmarking in horizontal analysis. We include six datasets spanning two major task domains—knowledge recall and reasoning. The number within parentheses denote the portion size of *non-answerable* queries given the provided context.

the *Qwen3 model family* (from 0.6B to 235B parameters; Yang et al., 2025) and *Deepseek-R1 distilled models* (DeepSeek-AI et al., 2025) based on Qwen2.5 (Yang et al., 2024) and Llama-3 (Dubey et al., 2024) backbones (from 1.5B to 70B parameters). For the Qwen3 models, we generate responses in both *thinking* and *non-thinking* modes to create a controlled comparison. See the full list of backbone models at Section A.

Our experiments focus on **simple queries** for which LLMs are more likely to exhibit overthinking (Shojaee et al., 2025). These queries span two main domains, sourced from six datasets (dataset statistics are summarized in Table 1).

- **Simple Reasoning:** Tasks requiring basic logical steps, such as grade-school math word problems (Miao et al., 2020), date arithmetic (Tan et al., 2023), and logic grid puzzles (Lin et al., 2025).
- **Knowledge Recall:** Tasks that necessitate recalling information with minimal reasoning: discerning unanswerable questions (Rajpurkar et al., 2018), locating facts from long contexts (Kamradt, 2023), and producing factually correct responses (Wei et al., 2024).

To ensure deterministic and reproducible outputs, we employ a **greedy decoding** strategy with temperature set to 0 and top-p set to 0. The maximum token generation limit was set to $8k$ tokens for thinking LLMs (or $16k$ for more complex temporal reasoning tasks) to allow for exhaustive exploration of the thought process, and $2k$ for the non-thinking mode.

Preliminary Results. Benchmark performances are shown in Figure 1 and Table A2 (full results). For simple reasoning tasks, the performance gains from enabling a model’s thinking capabilities are most significant for models with fewer than 4–8 billion parameters (Figure 1). Beyond this size threshold, the performance gap between thinking and non-thinking modes diminishes, mostly approaching zero. This demonstrates that additional thinking ceases to be effective and thus renders overthinking once the model size is beyond a threshold. Con-

versely, for knowledge recall tasks that involve a minimal reasoning workload, enabling thinking provides negligible benefits, irrespective of task difficulty (Figure A3 to A5).

Although these analyses focus on simple queries, the results already suggest that the necessity for long-form thinking is not dictated by the *overall task complexity*.⁵ Rather, it depends on the *expected reasoning workload*,⁶ which is further confirmed in Section 3.2 by varying reasoning efforts.

3.2 Vertical Analysis

Our *vertical* analysis narrows the focus to two domains (mathematical and temporal reasoning) where task difficulty can be systematically controlled. The aim here is to *identify the range of which extended reasoning might have become beneficial, and to distinguish it from settings where it offers little to no advantage*.

Experimental Setup. We adopt the same evaluation setup as in Section 3.1, testing 14 models with greedy decoding. For data domains, we focus on mathematical and temporal reasoning. Specifically, for mathematical reasoning, we use ASDiv (Miao et al., 2020) at grade 1–5 levels, along with GSM8k (Cobbe et al., 2021), which is considered more challenging than grade-5 ASDiv and solvable by bright middle-school students. For temporal reasoning, the original data provided in Tan et al. (2023) corresponds to level-1 difficulty. We then progressively increase query complexity using the procedure described in Section D, resulting in five difficulty levels. Sample examples for each level are shown in Table A1.

Preliminary Results. Our vertical study shows that *extended reasoning only pays off within narrow boundaries*—outside of which it quickly devolves into overthinking.

⁵The literature commonly uses task complexity as a signal to decide the expected reasoning efforts (Pu et al., 2025), including dimensions like reasoning workload and the amount of knowledge that needs to be recalled.

⁶Reasoning workload refers to amount of reasoning efforts, e.g., intermediate stepwise inference, required to reach a correct answer.

Difficulty Level	Performance	Efficiency	Tokens Waste
ASDiv-1	97.44/100.00	46.5/320.0 (255.8)	249.4
ASDiv-2	92.94/95.88	54.4/348.5 (274.5)	266.1
ASDiv-3	89.00/96.75	59.7/429.8 (349.7)	321.7
ASDiv-4	83.39/93.36	74.8/553.1 (450.8)	402.6
ASDiv-5	76.03/90.41	127.3/701.5 (553.9)	465.5
GSM8k	74.75/91.50	118.1/1,021.7 (889.7)	726.8

Table 2: Accuracy performance and inference efficiency (unit: words) of Qwen3-235B-A22B for vertical analysis (math reasoning). Values are reported as non-thinking/thinking, with thought lengths shown within parentheses. Full results see Table A3. Additionally, we introduce a summary metric, *token waste*, defined as the thinking tokens count multiplied by the performance ratio. Although this metric is not normalized, a larger value generally reflects a more severe token waste issue.

For *mathematical reasoning* (Table 2), thinking improves accuracy as tasks get harder and maintains its performance above 90%: from negligible gains at ASDiv-1/2 to 15 points at GSM8k. Yet this comes at a steep cost: solving GSM8k requires over $10\times$ more thought tokens, and still 80% of that extra compute produces no measurable gain. In other words, *even where thinking helps, the majority (e.g., 80%) of computation is wasted*. For *temporal reasoning* (Table 3), at L1-L2, where queries involve manageable duration representations, non-thinking models already achieve near-perfect accuracy, so extra reasoning adds little help. Beyond L3, where tasks demand day-level counting over hundreds or thousands of days, robust handling of leap years, and potential confusion with Julian calendar system, thinking performance collapses despite *huge reasoning workload*. While enabling thinking does help, boosting accuracy up to roughly 50%, the improvement quickly saturates. This reflects a ceiling imposed by the model’s **representational capacity**—its ability to internally encode and manipulate the structural rules required for the task. Past this capacity, *more thinking becomes pure overthinking and additional reasoning cannot bridge the gap*.

In summary, thinking is valuable only in a *narrow middle ground*: trivial tasks waste it, and tasks beyond the model’s representational capacity nullify it. While math domain highlights the first case, non-math domains such as temporal reasoning exhibits the second, revealing the true landscape of overthinking that math domain-intense studies alone cannot capture (i.a., Chen et al., 2024).

4 Analysis Framework of Inner Workings

Figure 2 introduces the overview of our proposed framework, TRACE, to study the inner work-

Difficulty Level	Performance	Efficiency
Temporal-L1	95.87/99.41	71.2/261.0 (162.0)
Temporal-L2	95.52/97.61	114.2/679.7 (534.7)
Temporal-L3	32.24/52.54	309.2/2,485.6 (2,255.8)
Temporal-L4	14.63/47.76	334.6/2,843.8 (2,625.4)
Temporal-L5	8.66/45.37	323.9/3,336.0 (3,078.9)

Table 3: Accuracy performance and inference efficiency (unit: words) of Qwen3-235B-A22B for vertical analysis (temporal reasoning). Full results see Table A4.

ings of an LLM’s thought process. TRACE, a **Thought-process Reconstruction and Automated Clustering Engine**, contains four main stages: (1) *Response Sampling*, where we generate responses in auto-regressive manner from various models; (2) *Thought Decomposition & Label Inference*, where we break down the reasoning and infer the relationships between sub-thoughts; (3) *Progression Graph Construction*, where we represent the reasoning flow in a graph structure; and (4) *Thought Pattern Induction*, where we aggregate individual graphs to discover generalized thought patterns.

For stage 1, we follow the same setting as in Section 3.1 to sample responses from 4 large thinking models: Qwen3-30B-A3B, Qwen3-32B, R1-Distill-Llama-70B and Qwen3-235B-A22B.⁷

Stage 2: Thought Decomposition & Label Inference. Once responses are collected, we utilize gemini-2.5-pro (Comanici et al., 2025) to systematically decompose each thought process into sequential **sub-thoughts** and infer the functional relationship between them. Specifically, we define a sub-thought as a text segment satisfying three strict criteria—self-contained, complete and answer-bearing (detailed in Section F.1)—and leverage sub-thought transition markers (Hamoud et al., 2025) to help decide the sub-thought boundaries.

Label Inference: While prior literature has informally described model behaviors (Gandhi et al., 2025), these descriptions often lack formal definitions. We establish a comprehensive and formally defined set of thought-to-thought relational labels to systematically categorize the transitions in the thought process. We have unary labels including *Initial* (first attempt, planning) and *Final* (final answer delivery); For the intermediate sub-thoughts, we devise binary labels: *Verification* (confirm the preceding thought’s correctness), *Correction* (rectify the preceding thought), *Backtrack* (reverts to an

⁷We focus on large models ($> 20B$) because smaller ones are unable to produce meaningful outputs for complex tasks. R1-Distill-Qwen-32B is not suited because it lacks the long-thinking ability required to explore thinking structure.

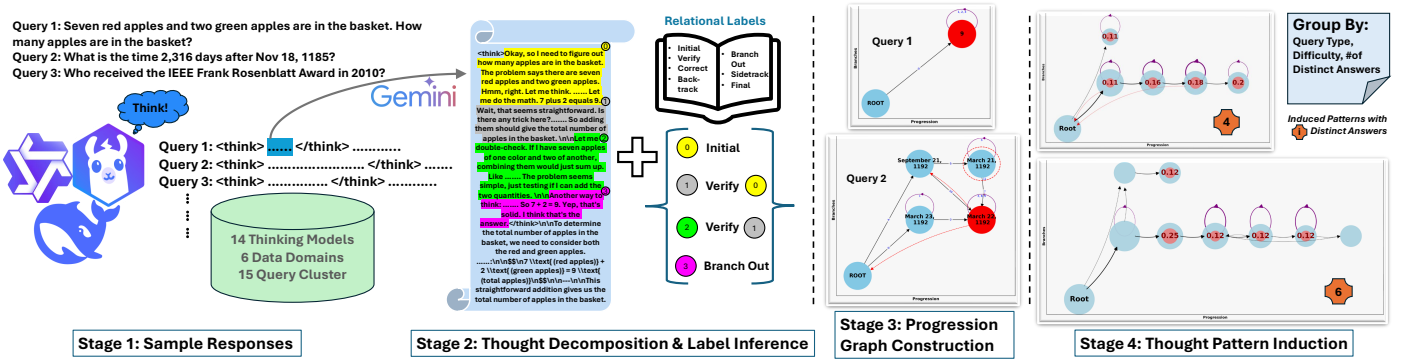


Figure 2: Overview of our proposed analyzer (TRACE) to study the inner workings of an LLM’s thought process. It contains four main stages (detailed in Section 4): *Response Sampling*, *Thought Decomposition & Label Inference*, *Progression Graph Construction*, and *Thought Pattern Induction*.

earlier path) and *Branching Out* (explore a new approach). *Sidetrack* is a unique label for rambling—a digressive, seemingly unrelated tangent that often adds no value to the thinking. Detailed definitions of each label are included in Section F.2.

We have provided the complete instruction in Figure A32 and A33 for reproducibility purpose. In addition to Figure 2 (stage 2), we provide another parsed output in Figure A34 to A35. Additionally, we also perform human inspection on 200 randomly sampled sub-thoughts, and we find that the automatic labels reasonable in 93% of the time.

Stage 3: Progression Graph Construction. To formally represent the reasoning trajectory, we construct a **thought progression graph** (e.g., Figure 3) for each individual thinking trace. In this graph, each node corresponds to a distinct answer proposed by the LLM during its thought process. A directed edge connects two nodes, representing the inferred relational label (e.g., *verification*) between the corresponding sub-thoughts. To facilitate systematic analysis, we project these graphs onto a 2D coordinate system, with a designated root node at (0,0) representing the query. Unary labels are encoded as follows: *initial* is marked as a directed edge from the root to the first explored answer at (1,1), and the *final* node (if applicable) is highlighted with a red dashed circle. For visualization, *branching out* is converted into either *verification* or *correction*, depending on whether the newly explored method produces the same or a different answer from the preceding node.

Regarding the layout, each new distinct answer is placed one step to the right (x -value + 1). When *backtracking* occurs and leads to exploration of a new answer, the new node is instead placed on an

elevated row (y -value + 1). This design choice serves two purposes: (1) it visually distinguishes the new, alternative path from prior linear progression, preventing overlaps, and (2) it emphasizes the **search depth** of the thinking process, with higher rows representing deeper layers of exploration prompted by prior *abandoned* path.

Stage 4: Thought Pattern Induction (Clustering). The final stage of TRACE aims to induce generalized and interpretable reasoning patterns (e.g., Figure 4) for a collection of topically similar individual progression graphs (Chambers and Jurafsky, 2008; Jin et al., 2022; Zhang et al., 2024).⁸ This is accomplished via a three-step process:

1. **Group-Based Aggregation:** Graphs are first grouped by shared characteristics, i.e., *query type*, *query difficulty*, and *the number of distinct answers proposed in the thought process*. Within each group, we perform a trivial alignment of the graphs using their coordinate representations and aggregate them. During aggregation, node weights are determined by their presence count, and edge weights are determined by their frequency.
2. **Thresholding:** We apply a heuristic threshold (e.g., 0.3) to the aggregated graph to prune infrequent nodes and edges. This filtering step removes noise and highlights the dominant reasoning paths.
3. **Ground-Truth Analysis:** Finally, we link the resulting patterns with the ground-truth answers. This allows us to analyze the distribution of ground-truth answers across different reasoning structures (e.g., *Late Landing* vs.

⁸Following the literature, we group responses by similar prompts, not responses to the identical prompt.

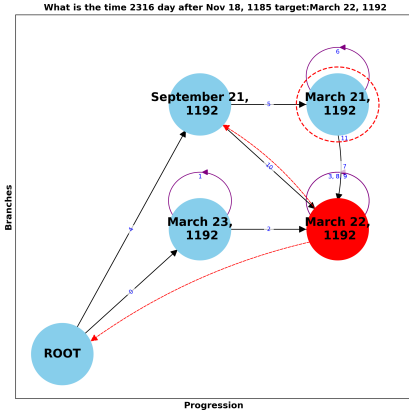


Figure 3: Individual thought progression graph of Qwen3-235B-A22B model on a sampled date arithmetic (temporal-L3) query. **Red bubble** denotes the ground-truth answer, while the **red dashed circle** denotes the final delivered answer.

Explorer; Section 5.1) and understand how different models arrive at their conclusions.

5 Findings

By applying our TRACE, we can systematically break down a model’s thought process into a quantifiable progression. Figure 3 presents one individual progression graph on a sample response.⁹ We then conduct a comprehensive analysis based on the *aggregation instead of individuals*, as the latter may display noisy or mixed behaviors. More importantly, our analysis reveals two predominant patterns of thought progression, when the model generates *at least 3 intermediate answers while thinking*.¹⁰ These patterns, which we term **Explorer** and **Late Landing**, reflect *intrinsic dynamics of models*—rather than properties of individual prompts—and represent different manifestations of their reasoning and tendencies to overthink.

Novel Definition of Overthinking Based on the thought patterns identified in Section 5.1 and the associated utility tracing (detailed in Section 5.2), we provide a refined definition for overthinking based on the decomposed *thought structure* (Figure 3).

Structure-based Overthinking Definition

*Overthinking is the continuation of thought beyond the point where the marginal **return** ($\Delta\text{Performance} / \Delta\text{Thought}$) drops below a predefined threshold ϵ .*

⁹Additional graphs are presented in Section G.

¹⁰For cases with ≤ 2 distinct answers, the progression is typically linear without branches. See Figure A17 and A23.

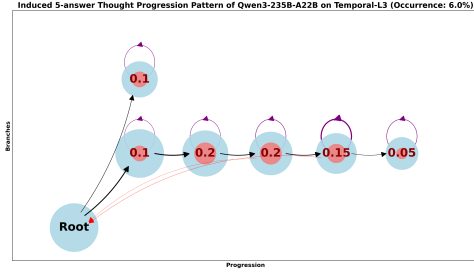


Figure 4: The typical **Explorer** thought progression pattern (5 distinct answer case). The size of the blue nodes indicates the visit frequency, while the size of the red nodes (and associated values) indicates the probability of the ground truth being present at that node. Due to the exploratory behavior, *multiple reasoning branches emerge and the correct answer can be discovered at any stage of the thought process*. Edge thickness indicates the edge frequency, and **red dashed curve** denotes the occurrence of backtracking, where the model abandon its current reasoning path. More plots in Section H.1.

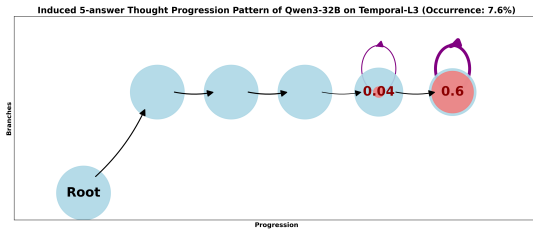


Figure 5: The typical **Late Landing** thought progression pattern (5 distinct answer case). The model follows a more linear path, with the probability of the ground-truth answer (indicated by the red node size and value) being highly concentrated at the terminal stage of the thought process. Towards the end, *the model engages in over-verification, marked by a thick self-loop, to increase its confidence*. More plots in Section H.2.

This redefinition formally identifies overthinking as reasoning that continues after the onset of *significant diminishing returns*, a critical threshold that we term the **convergence point**. Typically, $\Delta\text{Thought}$ corresponds to a one-sub-thought increment. Meanwhile, Section 5.3 details approximating the *return* via graph-induced characteristics.

5.1 Thought Progression Patterns

The two identified patterns are distinguished by *how the probability of arriving at the correct answer is distributed* throughout the thought process.

Explorer. This pattern (Figure 4) is characterized by the ground-truth answer being spread out across almost all nodes (i.e., distinct answers) in the thought process. In this mode, nearly every distinct intermediate answer has a non-trivial probability of being correct. A potential side effect of this ex-

Tracing Utility between Explorer (Left) and Late Landing (Right) Thinking Patterns

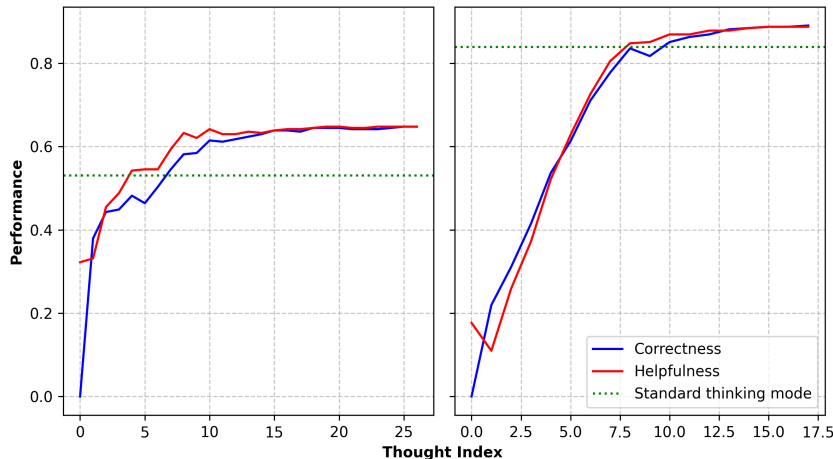


Figure 6: Tracing utilities on *Temporal-L3* task between **Explorer** (Qwen3-235B-A22B) and **Late Landing** (e.g., Qwen3-30B) thinking patterns. For **Explorer**, performance is volatile and peaks early, showing that further reasoning provides diminishing or even negative returns; while for **Late Landing**, it rises steadily to a plateau, and overthinking manifests as redundant steps after this convergence point. *Thinking mode* denotes the model performance when no intervention is applied.

ploratory behavior is that a correct answer might be discovered early in the process. As a result, this pattern of overthinking manifests as *excessive exploration*, where the model evaluates numerous alternatives, often leading it to refute its earlier conclusions, a behavior dubbed *backtrack*. While this exploration can cause the model to discard correct answers prematurely, returning to an earlier explored answer after backtracking significantly boosts the answer’s credibility, which is analogous to reaching the same answer via two independent, parallel reasoning paths.

Qwen3-235B-A22B, the largest evaluated model, exhibits this *exploratory* pattern. Note, the pattern is tied to models, irrespective of input prompts in general (see Section H.1 for plots on more tasks).

Late Landing. This pattern follows a “convergent” reasoning trajectory. As the thought process unfolds, the model gets closer to the correct answer, and the correctness probability increasingly concentrates towards the final proposal. As evidenced in Figure 5, early answers have extremely low correctness probabilities, while the final node overwhelmingly captures the correct answer. This behavior reflects a process of *sequential-esque correction*, in contrast to the branching exploration observed in the Explorer pattern. Overthinking in this pattern manifests as *excessive usage of self-verification*.¹¹ the model engages in an unnecessarily long chain of validating an already-correct answer, aiming for

overly cautious confidence rather than stopping once the confidence is adequate.

Most open-weight thinking models, such as R1-Distill-Llama-70B, Qwen3-30B-A3B, and Qwen3-32B, exhibit this *over-verification* pattern.

5.2 Utility Tracing

To quantify overthinking more structurally, we look into *utility tracing* by investigating the evolving performance as a function of the number of sub-thoughts. This analysis reveals distinct utility curves for our two previously identified thought patterns. Note, each sub-thought is analyzed for both **correctness** (whether the sub-thought contains the correct answer itself) and **helpfulness** (whether the thoughts can lead to a correct answer in the *final response*, i.e., post-`</think>` part shown in Figure 2).¹²

As shown in Figure 6, the **Explorer** pattern displays volatile performance that often peaks early, which means excessively extended reasoning can yield diminishing or negative returns. Conversely, the **Late Landing** pattern shows a steady initial performance increase (i.e., consistent slope); this progression suddenly converges to a stable plateau after a substantial number of sub-thoughts, where overthinking manifests as redundant steps (e.g., unnecessary verifications) past the *convergence point*.

¹¹The *Explorer* also exhibits over-usage of verification, but second to over-exploration.

¹²To evaluate helpfulness, we wrap prior sub-thoughts in `<think>` tokens, prompt the model to output responses only.

Model	Qwen3-235B-A22B		Qwen3-32B	
	Performance	Efficiency	Performance	Efficiency
No Thinking	32.33	263	17.68	427
Standard Thinking	52.87	2,722	83.84	4,000
Heuristic (SL k=2)	62.24	1,315	68.90	1,874
Heuristic (SL k=3)	63.44	1,681	80.18	2,463
Heuristic (SL k=2 + BT)	62.24	1,100	62.80	1,709

Table 4: Inference-time overthinking management using heuristics (SL: self-looping, BT: backtrack) based on our structure-based definition. Performance (accuracy) and efficiency (generation length) reveal that the optimal heuristic varies across models and is contingent on the thought-progression patterns.

5.3 Towards Overthinking Management

We present a case study on the *Temporal-L3* task to illustrate how our structure-based definition can be used to detect and manage overthinking in practice.

Based on the utility evolution of Figure 6, we empirically set $\epsilon = 0$ for Qwen3-235B-A22B and $\epsilon = 1$ for Qwen3-32B.¹³ For both models, the convergence point is the completion of the *eighth sub-thought*, classifying any additional thinking beyond this point as overthinking. Concretely, Figure 6 shows that the introduction of an additional thought reduces performance (63.25 \rightarrow 62.05) or provides negligible benefit (84.76 \rightarrow 85.06). Notably, performances at this convergence point already surpasses the standard thinking baselines (52.87 and 83.84, respectively).¹⁴ This highlights that our redefinition detects convergence point, and thus identifies overthinking at a fine-grained, sub-thought level, offering a more robust measurement of overthinking than the prior sample-level, length-based definition. In fact, the prior definition can be viewed as a special case of our novel redefinition, where $\Delta Thought$ is the *difference in total length* between thinking and non-thinking outputs.

Further, to extend this redefinition to real-world inference settings without ground-truths, we introduce two heuristics derived from the behavioral patterns observed in Section 5.1. They serve as proxies for estimating the convergence point.

- **Self-looping (SL)**: terminate once the model performs k consecutive self-verifications after proposing an answer (i.e., visit a node).
- **Backtrack (BT)**: terminate if and only if the model revisits a previously proposed answer *as the direct result of a backtrack action*.

As displayed in Table 4, we start with using self-looping heuristic only ($K = 2$), and we achieve

¹³We ignore the initial descent for one sub-thought case.

¹⁴4 samples were excluded due to gemini-2.5-pro parsing errors, accounting for minor discrepancies with Table 3.

62.23 and 68.90 accuracies for Qwen3-235B-A22B and Qwen3-32B, respectively, while average output length is halved (from 2,722 \rightarrow 1,315 and 4,000 \rightarrow 1,874). Model-specific tuning yields further improvements. For the Explorer-type Qwen3-235B-A22B, adding backtracking preserves accuracy while reducing length to 1,100 words (nearly 60% efficiency savings). For the Late-Landing-type Qwen3-32B, setting $k = 3$ slightly encourages verification, raising accuracy to 80.18 (3 points below the thinking mode) while cutting the inference cost by 40% ($4k \rightarrow 2,463$).

Overall, this case study showcases that our structure-based redefinition not only pinpoints overthinking in post-hoc analysis, but also enables practical heuristics for managing overthinking at real-time inference without access to ground truths.

6 Conclusion

In this work, we present TRACE, a fine-grained analyzer that reconstructs LLM reasoning into sub-thoughts and progression graphs, uncovering the structural evolution of (over)-thinking. Our large-scale preliminary study shows that models waste 5–20 \times more compute on simple tasks without accuracy gains, driven primarily by *over-verification* and *over-exploration*. Based on these insights, we propose a utility-based redefinition of overthinking, grounded in thought structures that moves beyond length-based metrics. Together, our TRACE and findings provide a sharper lens into LLMs’ internal decision-making and overthinking management.

Limitation

GPU resources. The LLMs utilized in this work are of 0.6 to 235 billion parameters. It is thus more time-consuming, especially the 235-billion model, than evaluating smaller previous-generation models like BART (Lewis et al., 2020), which in turn results in a significantly higher carbon footprint. Specifically, we evaluate models on NVIDIA H100 (80GB VRAM) or A100 (80GB VRAM), both with significant CPU and memory resources. For example, running inference with the Qwen3-235B-A22B model requires 8 NVIDIA H100s. The evaluation time for each model ranges from several minutes to 1 day, depending on the configurations (e.g., dataset, model size, thinking mode).

Evaluation Domains. In this work, we have included six evaluation domains, aiming to cover a diverse array of tasks and input–output mappings.

Yet, these domains do not comprehensively represent the full spectrum of NLP tasks. Indeed, evaluation remains an ongoing challenge for analyzing thinking models—existing studies have typically restricted their focus to the math domain only (Chen et al., 2024; Hammoud et al., 2025).

While our current evaluation analysis primarily involve relatively simple, well-posed prompts where extended reasoning may not yield clear advantages, future work will expand TRACE to more complex settings where creative or multi-step reasoning can provide tangible benefits. This broader coverage will help assess whether the observed thinking dynamics generalize to more challenging tasks (e.g., competitive coding) without conflating productive reasoning with overthinking.

References

- Pranjal Aggarwal, Seungone Kim, Jack Lanchantin, Sean Welleck, Jason Weston, Ilya Kulikov, and Swarnadeep Saha. 2025. *Optimalthinkingbench: Evaluating over and underthinking in llms*. *arXiv preprint arXiv:2508.13141*.
- Pranjal Aggarwal and Sean Welleck. 2025. *L1: controlling how long A reasoning model thinks with reinforcement learning*. *CoRR*, abs/2503.04697.
- Daman Arora and Andrea Zanette. 2025. *Training language models to reason efficiently*. *CoRR*, abs/2502.04463.
- Nathanael Chambers and Dan Jurafsky. 2008. *Unsupervised learning of narrative event chains*. In *Proceedings of ACL-08: HLT*, pages 789–797, Columbus, Ohio. Association for Computational Linguistics.
- Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. 2025. *Towards reasoning era: A survey of long chain-of-thought for reasoning large language models*. *CoRR*, abs/2503.09567.
- Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2024. *Do NOT think that much for 2+3=? on the overthinking of o1-like llms*. *CoRR*, abs/2412.21187.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. *Training verifiers to solve math word problems*. *CoRR*, abs/2110.14168.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit S. Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-Jiang Jiang, and 81 others. 2025. *Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities*. *CoRR*, abs/2507.06261.
- Sergio Da Silva. 2023. *System 1 vs. system 2 thinking*. *Psych*, 5(4):1057–1076.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 81 others. 2025. *Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning*. *CoRR*, abs/2501.12948.
- Yuntian Deng, Yejin Choi, and Stuart M. Shieber. 2024. *From explicit cot to implicit cot: Learning to internalize cot step by step*. *CoRR*, abs/2405.14838.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, and 82 others. 2024. *The llama 3 herd of models*. *CoRR*, abs/2407.21783.
- Gongfan Fang, Xinyin Ma, and Xinchao Wang. 2025. *Thinkless: LLM learns when to think*. *CoRR*, abs/2505.13379.
- Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh, Nathan Lile, and Noah D. Goodman. 2025. *Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars*. *CoRR*, abs/2503.01307.
- Hasan Abed Al Kader Hammoud, Hani Itani, and Bernard Ghanem. 2025. *Beyond the last answer: Your reasoning trace uncovers more than you think*. *CoRR*, abs/2504.20708.
- Tingxu Han, Zhenting Wang, Chunrong Fang, Shiyu Zhao, Shiqing Ma, and Zhenyu Chen. 2025. *Token-budget-aware LLM reasoning*. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 24842–24855, Vienna, Austria. Association for Computational Linguistics.
- Michael Hassid, Gabriel Synnaeve, Yossi Adi, and Roy Schwartz. 2025. *Don't overthink it. preferring shorter thinking chains for improved LLM reasoning*. *CoRR*, abs/2505.17813.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. *Measuring mathematical problem solving with the MATH dataset*. In *Proceedings of the Neural Information Processing Systems Track*

- on Datasets and Benchmarks 1, *NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helvar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Ifimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, and 80 others. 2024. [Openai o1 system card](#). *CoRR*, abs/2412.16720.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2025. [Livecodebench: Holistic and contamination free evaluation of large language models for code](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- Mingyu Jin, Qinkai Yu, Dong Shu, Haiyan Zhao, Wenyue Hua, Yanda Meng, Yongfeng Zhang, and Mengnan Du. 2024. [The impact of reasoning step length on large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 1830–1842, Bangkok, Thailand. Association for Computational Linguistics.
- Xiaomeng Jin, Manling Li, and Heng Ji. 2022. [Event schema induction with double graph autoencoders](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2013–2025, Seattle, United States. Association for Computational Linguistics.
- Gregory Kamradt. 2023. [Needle in a haystack – pressure testing llms](#). https://github.com/gkamradt/LLMTest_NeedleInAHaystack. GitHub repository.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. [Large language models are zero-shot reasoners](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. [Efficient memory management for large language model serving with pagedattention](#). In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Zhong-Zhi Li, Duzhen Zhang, Ming-Liang Zhang, Jiaxin Zhang, Zengyan Liu, Yuxuan Yao, Haotian Xu, Junhao Zheng, Pei-Jie Wang, Xiuyi Chen, Yingying Zhang, Fei Yin, Jiahua Dong, Zhijiang Guo, Le Song, and Cheng-Lin Liu. 2025. [From system 1 to system 2: A survey of reasoning large language models](#). *CoRR*, abs/2502.17419.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. [Let’s verify step by step](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Bill Yuchen Lin, Ronan Le Bras, Kyle Richardson, Ashish Sabharwal, Radha Poovendran, Peter Clark, and Yejin Choi. 2025. [ZebraLogic: On the scaling limits of llms for logical reasoning](#). *CoRR*, abs/2502.01100.
- Ryan Liu, Jiayi Geng, Addison J. Wu, Ilia Sucholutsky, Tania Lombrozo, and Thomas L. Griffiths. 2024. [Mind your step \(by step\): Chain-of-thought can reduce performance on tasks where thinking makes humans worse](#). *CoRR*, abs/2410.21333.
- Thang Luong and Edward Lockhart. 2025. [Advanced version of gemini with deep think officially achieves gold-medal standard at the international mathematical olympiad](#). Blog post.
- Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su. 2020. [A diverse corpus for evaluating and developing English math word problem solvers](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 975–984, Online. Association for Computational Linguistics.
- Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel J. Candès, and Tatsunori Hashimoto. 2025. [s1: Simple test-time scaling](#). *CoRR*, abs/2501.19393.
- Alhassan Mumuni and Fuseini Mumuni. 2025. [Large language models for artificial general intelligence \(AGI\): A survey of foundational principles and approaches](#). *CoRR*, abs/2501.03151.
- Xiao Pu, Michael Saxon, Wenyue Hua, and William Yang Wang. 2025. [THOUGHTTERMINATOR: benchmarking, calibrating, and mitigating overthinking in reasoning models](#). *CoRR*, abs/2504.13367.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. [Know what you don’t know: Unanswerable questions for SQuAD](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia. Association for Computational Linguistics.

- Stephen E. Robertson and Hugo Zaragoza. 2009. [The probabilistic relevance framework: BM25 and beyond](#). *Found. Trends Inf. Retr.*, 3(4):333–389.
- Si Shen, Fei Huang, Zhixiao Zhao, Chang Liu, Tian-sheng Zheng, and Danhao Zhu. 2025. [Long is more important than difficult for training reasoning models](#). *CoRR*, abs/2503.18069.
- Quan Shi, Michael Tang, Karthik Narasimhan, and Shunyu Yao. 2024. [Can language models solve olympiad programming?](#) *CoRR*, abs/2404.10952.
- Parshin Shojaee, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. 2025. [The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity](#). *CoRR*, abs/2506.06941.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. [Scaling LLM test-time compute optimally can be more effective than scaling model parameters](#). *CoRR*, abs/2408.03314.
- Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Shaochen Zhong, Hanjie Chen, and Xia Ben Hu. 2025. [Stop overthinking: A survey on efficient reasoning for large language models](#). *CoRR*, abs/2503.16419.
- Qingyu Tan, Hwee Tou Ng, and Lidong Bing. 2023. [Towards benchmarking and improving the temporal reasoning capability of large language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14820–14835, Toronto, Canada. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.
- Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese, John Schulman, and William Fedus. 2024. [Measuring short-form factuality in large language models](#). *CoRR*, abs/2411.04368.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Fengli Xu, Qian Yue Hao, Zefang Zong, Jingwei Wang, Yunke Zhang, Jingyi Wang, Xiaochong Lan, Jiahui Gong, Tianjian Ouyang, Fanjin Meng, Chenyang Shao, Yuwei Yan, Qinglong Yang, Yiwen Song, Sijian Ren, Xinyuan Hu, Yu Li, Jie Feng, Chen Gao, and Yong Li. 2025a. [Towards large reasoning models: A survey of reinforced reasoning with large language models](#). *CoRR*, abs/2501.09686.
- Silei Xu, Wenhao Xie, Lingxiao Zhao, and Pengcheng He. 2025b. [Chain of draft: Thinking faster by writing less](#). *CoRR*, abs/2502.18600.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chuji Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 40 others. 2025. [Qwen3 technical report](#). *CoRR*, abs/2505.09388.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 22 others. 2024. [Qwen2.5 technical report](#). *CoRR*, abs/2412.15115.
- Xinliang Frederick Zhang, Nick Beauchamp, and Lu Wang. 2024. [Narrative-of-thought: Improving temporal reasoning of large language models via re-constructed narratives](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 16507–16530, Miami, Florida, USA. Association for Computational Linguistics.
- Xinliang Frederick Zhang, Nick Beauchamp, and Lu Wang. 2025a. [PRIME: large language model personalization with cognitive memory and thought processes](#). *CoRR*, abs/2507.04607.
- Yunxiang Zhang, Muhammad Khalifa, Lechen Zhang, Xin Liu, Ayoung Lee, Xinliang Frederick Zhang, Farima Fatahi Bayat, and Lu Wang. 2025b. [Logit arithmetic elicits long reasoning capabilities without training](#). *CoRR*, abs/2507.12759.
- Chuji Zheng, Zhenru Zhang, Beichen Zhang, Runji Lin, Keming Lu, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. 2025. [Processbench: Identifying process errors in mathematical reasoning](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, *ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 1009–1024. Association for Computational Linguistics.

A Models Used in This Work

This work utilized a diverse set of large language models (LLMs), especially those with long thinking capabilities. The models were primarily chosen to cover a wide range of parameter sizes, from smaller, more efficient models to very large, high-performing ones.

The complete list of models used in this study is provided below:

- **Qwen3 models (supports both thinking and non-thinking modes):**

- Qwen3-0.6B
- Qwen3-1.7B
- Qwen3-4B
- Qwen3-8B
- Qwen3-14B
- Qwen3-30B-A3B
- Qwen3-32B
- Qwen3-235B-A22B

- **R1-distilled thinking models:**

- DeepSeek-R1-Distill-Owen-1.5B
- DeepSeek-R1-Distill-Owen-7B
- DeepSeek-R1-Distill-Owen-8B
- DeepSeek-R1-Distill-Owen-14B
- DeepSeek-R1-Distill-Owen-32B
- DeepSeek-R1-Distill-Llama-70B

We use vLLM (Kwon et al., 2023) to deploy aforementioned models.

B Data Used in This Work

This work utilizes 6 data domains, covering 15 query clusters, for systematic evaluations. The six domains and associated datasets used in the horizontal analysis are displayed in Table 1.

The complete list of data domains and query clusters being evaluated in this study is provided below:

- **Knowledge Recall:**

- Short-context Machine Reasoning Comprehension domain: SQuAD2.0 (Rajpurkar et al., 2018).¹⁵
- Long-context Fact Retrieval domain: NAIH (Kamradt, 2023).¹⁶
- Factuality Question Answering domain: SimpleQA (Wei et al., 2024).¹⁷

- **Reasoning:**

- Mathematical Reasoning domain (sorted by difficulty): ASDiv-grade1 (Miao et al., 2020), ASDiv-grade2 (Miao et al., 2020), ASDiv-grade3 (Miao et al., 2020), ASDiv-grade4 (Miao et al., 2020),

ASDiv-grade5 (Miao et al., 2020),¹⁸ and GSM8k (Cobbe et al., 2021).¹⁹

- Temporal Reasoning domain: Date arithmetic (a.k.a., Temporal-L1; Tan et al., 2023).²⁰ Temporal-L2 to L5 are derived from Tan et al. (2023) following our designed curation guideline (Section D).
- Logical Reasoning domain: Zebra Logic (Easy; Lin et al., 2025).²¹

Note, the usage of each dataset is consistent with its intended use, and these datasets present no more than minimal risk.

C More Plots for Horizontal Study

Figure A1, A2, A3, A4, and A5 display the performance and inference-time efficiency trends of Qwen3 models on temporal reasoning, logical reasoning, Short-context MRC, Long-context Fact Retrieval, and Factuality QA tasks, respectively.

D Multi-level Temporal Reasoning (Date Arithmetic) Data Curation Process

A date arithmetic question from the original temporal reasoning dataset (Tan et al., 2023) is like “What is the time 10 year and 1 month after Feb, 1068”, which contains two important components: starting time and duration.

We perform systematic adjustments as follows to gradually increase the task complexity:

- From level-1 to level-2, we add *day* information to both starting time and duration.
- From level-2 to level-3, we convert the duration from the standard representation (*year, month and day*) to day-only representation, i.e., 10 year, 0 month and 21 day is equivalent to 3,674 days if the starting date is Feb 26, 1068.
- For level-4 and level-5, we double and triple the duration, respectively. In order to determine target dates, we use online date calculator to ensure the process flawless.

We show sample examples in Table A1 for each difficulty level.

¹⁵https://huggingface.co/datasets/rajpurkar/squad_v2. The dataset is distributed under the CC BY-SA 4.0 license.

¹⁶https://github.com/gkamradt/LLMTest_NeedleInAHaystack. The dataset is distributed under the MIT license.

¹⁷<https://openai.com/index/introducing-simpleqa/>. The dataset is distributed under the MIT license.

¹⁸<https://github.com/chaochun/nlu-asdiv-dataset>. The dataset is distributed under the CC BY-NC 4.0 license.

¹⁹<https://huggingface.co/datasets/openai/gsm8k>. The dataset is distributed under the MIT license.

²⁰<https://huggingface.co/datasets/tonytan48/TempReason>. This dataset is released on a public GitHub.

²¹<https://huggingface.co/blog/yuchenlin/zebra-logic>. The dataset is distributed under the CC BY 4.0 license.

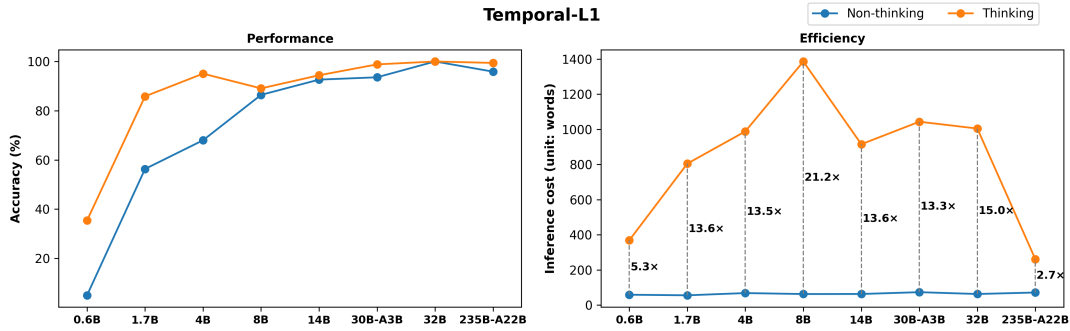


Figure A1: Performance and inference-time efficiency trends of Qwen3 models at different scales on simple temporal reasoning.

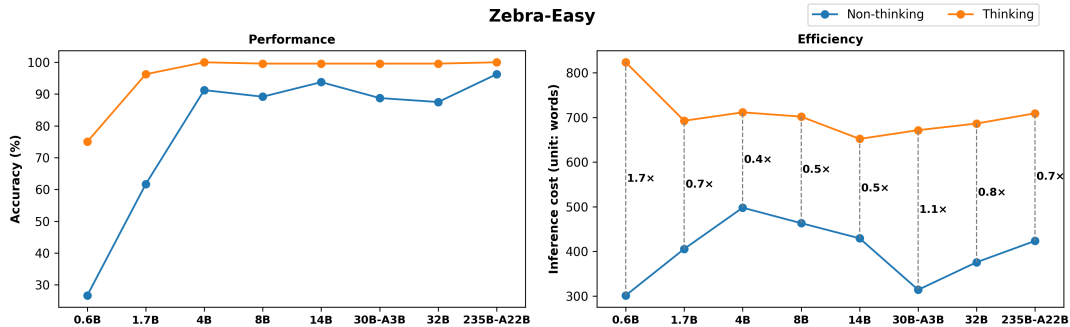


Figure A2: Performance and inference-time efficiency trends of Qwen3 models at different scales on simple logic reasoning.

E Full Benchmark Results of Section 3

Table A2, Table A3 and Table A4 display the full benchmark results using the Qwen3 model family.

Table A5, Table A6 and Table A7 display the full benchmark results using the DeepSeek-R1-Distill model family.

F Sub-thought Decomposition Details

F.1 Sub-thought Criteria

We define a sub-thought according to three strict criteria:

- **Self-contained:** It must be independently comprehensible without heavy reliance on other sub-thoughts.
- **Complete:** It must represent a full logical step, typically including an intent, an action or calculation, and a conclusion.
- **Answer-bearing:** It must explicitly state a proposed answer to the query.

F.2 Discourse Label Details

- **Initial:** The first attempt at solving the problem, which may include an outline of the overall plan. Each thought process has only one initial sub-thought.
- **Verification:** A confirmation of the correctness of the *immediately preceding* sub-thought, without any changes or corrections.

- **Correction:** A modification or rectification of the *immediately preceding* sub-thought.
- **Backtrack:** The abandonment of the current line of reasoning to revert to an earlier (but not immediately preceding) sub-thought and explore a different path.
- **Sidetrack:** An exploration of supplementary information that is not directly part of an alternative approach to the main question.
- **Branching Out:** An exploration of an alternative method or a different aspect of the problem, as opposed to merely adding information (sidetrack).
- **Final:** The sub-thought that delivers the final proposed answer.

G More Plots for Individual Thought Progression

This section presents selected individual thought progression plots.

Figure A6 displays the thought progression of the Qwen3-235B-A22B model on a grade-1 arithmetic query.

Figure A7 displays the thought progression of the R1-Distill-Llama-70B model on a grade-1 arithmetic query.

Figure A8 displays the thought progression of the Qwen3-32B model on a date arithmetic (temporal-L3) query.

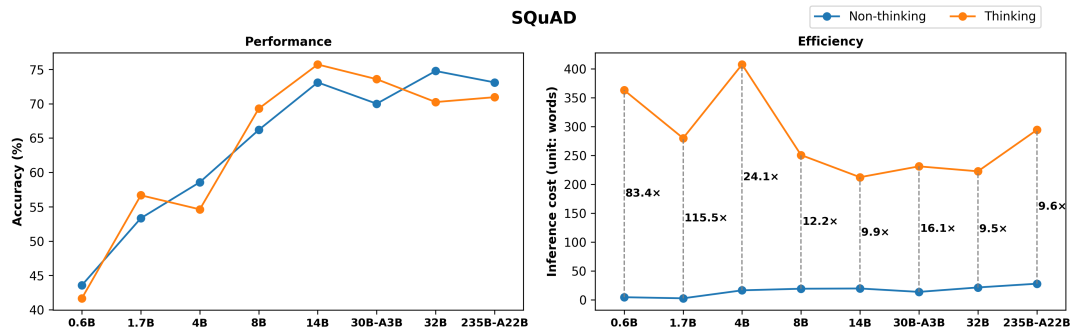


Figure A3: Performance and inference-time efficiency trends of Qwen3 models at different scales on simple short-context knowledge recall (machine reading comprehension).

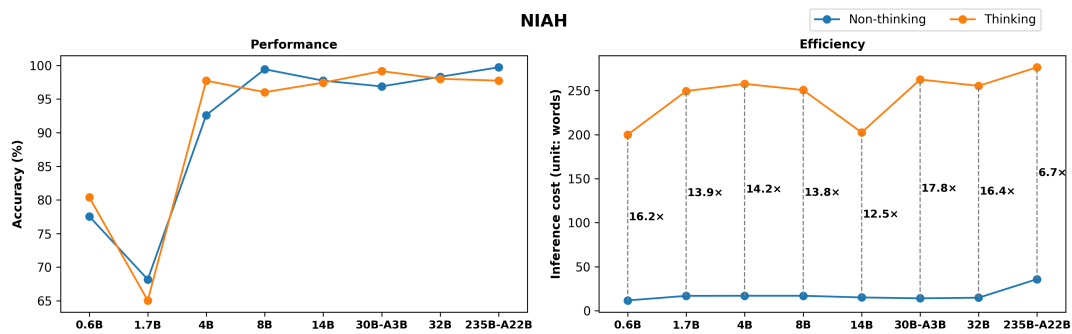


Figure A4: Performance and inference-time efficiency trends of Qwen3 models at different scales on simple long-context knowledge recall (fact retrieval).

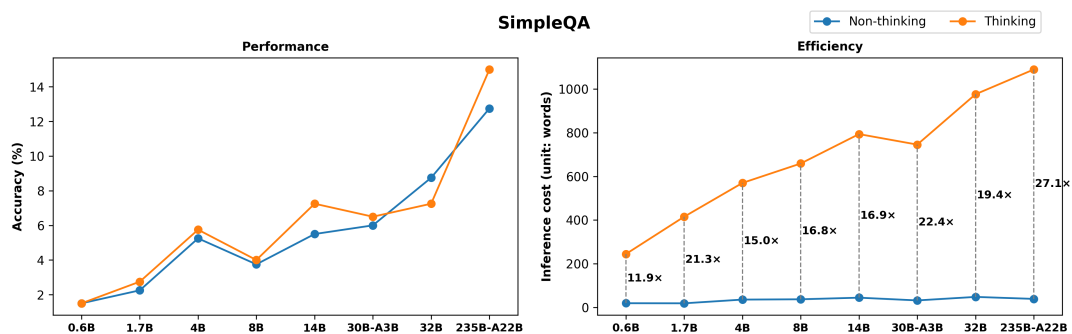


Figure A5: Performance and inference-time efficiency trends of Qwen3 models at different scales on a factuality QA task.

Difficulty Level	Question	Target
Level-1	What is the time 10 year and 1 month after Feb, 1068	Mar, 1078
Level-2	What is the time 10 year, 0 month and 21 day after Feb 26, 1068	Mar 19, 1078
Level-3	What is the time 3674 day after Feb 26, 1068	Mar 19, 1078
Level-4	What is the time 7348 day after Feb 26, 1068	Apr 9, 1088
Level-5	What is the time 11022 day after Feb 26, 1068	May 1, 1098

Table A1: Sample examples for different difficulty levels of the temporal reasoning task.

Model	SQuAD		NIAH		SimpleQA		ASDiv-1		Date Arithmetic (L1)		Zebra (Easy)	
	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency
Qwen3-0.6B	43.57/41.67	4.3/362.8(360.6)	77.56/80.40	11.6/200.0(192.4)	1.50/1.50	19.0/244.3(232.5)	87.69/97.44	39.1/174.1(148.8)	5.01/35.40	58.6/368.5(332.4)	26.67/75.00	301.3/823.6(642.6)
Qwen3-1.7B	53.33/56.67	2.4/279.7(275.9)	68.18/65.06	16.8/249.5(227.3)	2.25/2.75	18.6/415.3(384.9)	94.87/99.49	40.2/182.2(148.3)	56.34/85.84	55.3/804.8(720.3)	61.67/96.25	405.5/692.6(488.5)
Qwen3-4B	58.57/64.52	16.2/407.3(382.7)	92.61/97.73	16.9/257.7(229.0)	5.25/5.75	35.6/570.0(501.2)	99.49/100.00	44.5/222.5(178.6)	91.74/97.05	68.0/988.9(880.7)	91.25/100.00	497.9/711.5(486.2)
Qwen3-8B	66.19/69.29	19.0/250.5(226.6)	99.43/96.02	16.9/250.6(219.7)	3.75/4.00	37.0/659.0(562.9)	97.44/98.97	43.3/338.3(267.4)	86.43/89.97	62.5/1386.0(1265.5)	89.17/99.58	463.2/701.8(479.0)
Qwen3-14B	73.10/75.71	19.4/212.1(203.8)	97.73/97.44	15.0/202.5(180.5)	5.50/7.25	44.4/793.7(682.8)	96.92/98.46	41.8/274.1(229.8)	92.63/94.40	62.8/915.0(792.1)	93.75/99.58	429.2/652.0(428.6)
Qwen3-30B-A3B	70.00/73.57	13.5/230.9(204.1)	96.88/99.15	14.0/262.6(233.0)	6.00/6.50	31.8/745.6(644.5)	97.95/98.97	38.5/266.3(218.9)	98.53/98.82	73.1/1043.2(925.4)	88.75/99.58	313.9/671.5(451.4)
Qwen3-32B	74.76/70.24	21.2/222.3(201.3)	98.30/98.01	14.7/255.4(230.9)	8.75/7.25	47.9/976.2(797.0)	94.87/100.00	42.3/290.1(221.9)	100.00/100.00	62.7/1004.4(883.3)	87.50/99.58	375.7/686.3(459.6)
Qwen3-235B-A22B	73.10/70.95	27.7/294.0(254.8)	99.72/97.73	35.8/276.5(247.9)	12.75/15.00	38.8/1089.7(953.2)	97.44/100.00	46.5/320.0(255.8)	95.87/99.41	71.2/261.0(162.0)	96.25/100.00	423.5/907.3(678.1)

Table A2: Complete accuracy performance and inference efficiency (unit: words) of Qwen3 model family on horizontal analysis. For both performance and efficiency metrics, The first number is for the non-thinking mode, while the second is for the thinking mode. The numbers in parentheses indicate the thought lengths. Results of DeepSeek-R1-Distill see Table A5.

Figure A9 displays the thought progression of the Qwen3-32B model on a grade-1 arithmetic query.

Figure A10 displays the thought progression of the Qwen3-32B model on a GSM8k query.

Figure A11 displays the thought progression of the Qwen3-32B model on another GSM8k query.

Figure A12 displays the thought progression of the R1-Distill-Llama-70B model on a date arithmetic (temporal-L3) query.

Figure A13 displays the thought progression of the R1-Distill-Llama-70B model on a GSM8k query.

Figure A14 displays the thought progression of the Qwen3-235B-A22B model on a GSM8k query.

Figure A15 displays the thought progression of the Qwen3-235B-A22B model on another GSM8k query.

Figure A16 displays the thought progression of the Qwen3-235B-A22B model on a GSM8k query, which is the graphical representation of Figure A34 and A35.

H More Plots for Induced Thought Progression Patterns

H.1 More Plots of Explorer Thought Progression Pattern

Figure A17 (trivial case), A18, A19, A20, A21 (trivial case) and A22 present more plots of the **Explorer** thought progression pattern.

H.2 More Plots of Late Landing Thought Progression Pattern

Figure A23 (trivial case), A24, A25, A26 and A27 present more plots of the **Late Landing** thought progression pattern.

H.3 Comparison of Induced Thought Progression Patterns with Different Distinct Answers

Figure A28 to A31 show how the progression patterns change with respect to the increasing number of distinct answers. The selected task is *SimpleQA* which is factuality QA task with minimal reasoning efforts. Therefore, the extended reasoning is futile, and the model is trapped by over-verification and frequent revisiting of a previous answer (*marked by the extremely think back edge in Figure A29 and Figure A30*).

Model	ASDiv-1		ASDiv-2		ASDiv-3		ASDiv-4		ASDiv-5		GSM8k	
	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency
Qwen3-0.6B	87.69/97.44	39.1/174.1 (148.8)	86.18/91.76	45.0/242.3 (216.7)	77.75/90.25	45.7/274.2 (245.4)	70.43/80.73	53.0/312.0 (278.4)	55.48/82.19	110.0/412.0 (343.9)	50.75/74.75	84.9/582.4 (529.4)
Qwen3-1.7B	94.87/99.49	40.2/182.2 (148.3)	87.65/94.71	45.3/228.9 (192.4)	85.75/93.75	53.3/254.5 (212.9)	77.74/91.69	58.9/330.9 (275.8)	68.49/88.36	108.8/490.6 (392.8)	65.25/87.00	105.6/555.3 (472.7)
Qwen3-4B	99.49/100.00	44.5/222.5 (178.6)	95.00/95.88	51.4/286.6 (233.5)	91.00/96.50	55.7/346.9 (281.5)	83.06/95.35	69.4/454.4 (368.8)	77.40/90.41	121.4/695.9 (556.9)	76.00/91.75	112.0/786.1 (667.6)
Qwen3-8B	97.44/98.97	43.3/338.3 (267.4)	95.00/95.88	51.8/362.1 (282.8)	92.25/95.50	56.1/398.4 (310.5)	79.40/94.02	67.7/511.1 (404.0)	76.03/90.41	127.6/770.6 (588.5)	80.00/92.00	119.7/818.0 (681.6)
Qwen3-14B	96.92/98.46	41.8/274.1 (229.8)	93.82/96.18	50.3/313.3 (257.5)	88.00/96.25	52.8/351.8 (283.4)	81.06/97.01	71.3/499.2 (400.5)	73.29/90.41	134.3/829.1 (664.7)	74.25/93.50	110.8/781.0 (654.1)
Qwen3-30B-A3B	97.95/98.97	38.5/266.3 (218.9)	94.71/96.18	43.5/296.9 (239.8)	93.25/96.25	48.5/333.7 (265.7)	84.72/95.35	61.7/422.1 (330.8)	75.34/89.73	131.9/769.7 (604.6)	79.50/94.00	112.7/771.3 (639.8)
Qwen3-32B	94.87/100.00	42.3/290.1 (221.9)	90.59/96.18	52.1/346.7 (259.9)	84.50/96.25	58.4/394.1 (302.7)	78.74/94.35	77.3/515.4 (399.6)	68.49/86.99	170.2/903.9 (719.4)	65.25/92.25	118.2/764.2 (622.7)
Qwen3-235B-A22B	97.44/100.00	46.5/320.0 (255.8)	92.94/95.88	54.4/348.5 (274.5)	89.00/96.75	59.7/429.8 (349.7)	83.39/93.36	74.8/553.1 (450.8)	76.03/90.41	127.3/701.5 (553.9)	74.75/91.50	118.1/1021.7 (889.7)

Table A3: Complete accuracy performance and inference efficiency (unit: words) of Qwen3 model family on vertical analysis (math reasoning). For both performance and efficiency metrics, The first number is for the non-thinking mode, while the second is for the thinking mode. The number in parentheses indicates the thought length. Results of DeepSeek-R1-Distill see Table A6.

Model	Temporal-L1		Temporal-L2		Temporal		Temporal-L4		Temporal-L5	
	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency
Qwen3-4B	91.74/97.05	68.0/988.9 (880.7)	74.03/88.96	92.1/1748.1 (1608.6)	2.69/59.70	246.2/3933.3 (3708.7)	0.90/61.49	281.7/2200.7 (1944.7)	0.60/53.43	178.5/5134.5 (4878.0)
Qwen3-8B	86.43/89.97	62.5/1386.0 (1265.5)	77.31/90.75	84.2/2046.1 (1900.1)	8.96/81.79	352.9/4212.5 (3927.5)	3.58/80.60	335.2/3748.2 (3503.6)	2.39/82.69	372.0/3807.9 (3502.1)
Qwen3-14B	92.63/94.40	62.8/915.0 (792.1)	81.79/98.51	85.3/1924.3 (1739.0)	8.36/88.06	318.7/3408.5 (3161.4)	6.87/87.76	319.3/3823.8 (3568.2)	0.90/85.67	307.7/2993.3 (2706.3)
Qwen3-30B-A3B	98.53/98.82	73.1/1043.2 (925.4)	97.01/98.21	110.7/2306.3 (2147.4)	20.90/89.55	332.6/3684.5 (3421.2)	17.01/89.85	340.8/3624.8 (3369.5)	14.03/85.07	356.5/3601.8 (3334.9)
Qwen3-32B	100.0/100.00	62.7/1004.4 (883.3)	93.13/95.52	93.2/2198.1 (2014.5)	17.61/82.69	377.1/3992.6 (3714.7)	12.24/80.60	335.1/4001.9 (3738.7)	7.46/80.30	334.8/4217.6 (3963.2)
Qwen3-235B-A22B	95.87/99.41	71.2/261.0 (162.0)	95.52/97.61	114.2/679.7 (534.7)	32.24/52.54	309.2/2485.6 (2255.8)	14.63/47.76	334.6/2843.8 (2625.4)	8.66/45.37	323.9/3336.0 (3078.9)

Table A4: Complete accuracy performance and inference efficiency (unit: words) of Qwen3 model family on vertical analysis (temporal reasoning). For both performance and efficiency metrics, The first number is for the non-thinking mode, while the second is for the thinking mode. The number in parentheses indicates the thought length. Results of DeepSeek-R1-Distill see Table A7.

Model	SQuAD		NIAH		SimpleQA		ASDiv-1		Date Arithmetic (L1)		Zebra (Easy)	
	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency
DeepSeek-R1-Distill-Qwen-1.5B	35.48	317.0 (310.8)	5.11	2818.2 (2792.7)	0.00	- (-)	65.13	280.2 (253.4)	31.86	331.1 (315.8)	12.92	653.0 (470.7)
DeepSeek-R1-Distill-Qwen-7B	47.38	221.0 (210.6)	7.10	858.0 (835.7)	0.00	- (-)	93.85	119.6 (57.1)	33.92	146.2 (58.3)	69.58	855.6 (617.1)
DeepSeek-R1-Distill-Llama-8B	51.67	232.9 (224.9)	65.06	388.7 (371.5)	0.25	606.0 (583.0)	93.33	135.9 (64.7)	69.03	216.3 (70.9)	90.42	901.9 (667.7)
DeepSeek-R1-Distill-Qwen-14B	63.10	228.5 (207.1)	82.67	287.2 (262.9)	7.00	500.3 (480.9)	94.36	120.9 (59.2)	88.79	167.6 (65.3)	95.42	778.7 (562.0)
DeepSeek-R1-Distill-Qwen-32B	71.90	217.8 (198.1)	88.35	254.4 (234.3)	7.25	397.1 (377.8)	95.90	129.3 (61.3)	98.82	129.8 (60.1)	98.75	720.7 (497.6)
DeepSeek-R1-Distill-Llama-70B	69.29	235.6 (215.5)	93.75	213.6 (198.0)	20.00	401.6 (382.7)	95.90	121.4 (55.8)	93.22	243.9 (203.4)	98.75	733.1 (519.9)

Table A5: Complete accuracy performance and inference efficiency (unit: words) of DeepSeek-R1-Distill model family on horizontal analysis. The numbers in parentheses indicate the thought lengths.

Model	ASDiv-1		ASDiv-2		ASDiv-3		ASDiv-4		ASDiv-5		GSM8k	
	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency
DeepSeek-R1-Distill-Qwen-1.5B	65.13	280.2 (253.4)	49.12	278.6 (245.0)	50.75	298.0 (260.3)	33.22	321.6 (287.9)	39.04	392.1 (357.7)	27.00	372.2 (302.9)
DeepSeek-R1-Distill-Qwen-7B	93.85	119.6 (57.1)	91.47	134.1 (61.4)	81.25	146.6 (67.4)	77.74	170.2 (76.9)	75.34	236.2 (101.8)	72.25	239.1 (98.5)
DeepSeek-R1-Distill-Llama-8B	93.33	135.9 (64.7)	90.29	147.0 (68.4)	86.00	161.6 (76.4)	77.41	173.3 (80.7)	71.71	252.3 (114.4)	71.50	249.1 (111.9)
DeepSeek-R1-Distill-Qwen-14B	94.36	120.9 (59.2)	91.18	133.7 (62.7)	85.75	157.5 (72.1)	77.74	173.5 (78.4)	80.14	234.8 (101.2)	83.00	246.0 (106.6)
DeepSeek-R1-Distill-Qwen-32B	95.90	129.3 (61.3)	93.24	145.6 (70.1)	89.50	164.3 (73.5)	87.38	181.5 (80.9)	84.93	258.3 (115.4)	88.00	254.5 (106.2)
DeepSeek-R1-Distill-Llama-70B	95.90	121.4 (55.8)	92.65	136.6 (62.4)	86.50	154.2 (67.7)	84.39	174.0 (74.6)	74.66	252.3 (105.2)	79.50	263.7 (105.9)

Table A6: Complete accuracy performance and inference efficiency (unit: words) of DeepSeek-R1-Distill model family on vertical analysis (math reasoning). The numbers in parentheses indicate the thought lengths.

Model	Temporal-L1		Temporal-L2		Temporal-L3		Temporal-L4		Temporal-L5	
	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency	Performance	Efficiency
DeepSeek-R1-Distill-Qwen-1.5B	31.86	331.1 (315.8)	11.94	431.5 (413.3)	0.00	- (-)	0.00	- (-)	0.00	- (-)
DeepSeek-R1-Distill-Qwen-7B	33.92	146.2 (58.3)	67.76	182.9 (64.9)	2.39	453.9 (127.8)	1.79	1866.2 (155.0)	1.19	546.8 (125.0)
DeepSeek-R1-Distill-Llama-8B	69.03	216.3 (70.9)	59.10	265.0 (92.7)	1.49	419.4 (171.0)	0.30	332.0 (187.0)	1.19	315.5 (125.0)
DeepSeek-R1-Distill-Qwen-14B	88.79	167.6 (65.3)	55.22	209.8 (96.1)	4.18	486.0 (167.7)	3.58	424.3 (150.8)	2.99	402.0 (156.9)
DeepSeek-R1-Distill-Qwen-32B	98.82	129.8 (60.1)	91.34	153.2 (72.0)	3.28	444.1 (173.8)	4.78	509.3 (165.9)	2.69	444.7 (144.1)
DeepSeek-R1-Distill-Llama-70B	93.22	243.9 (203.4)	91.04	381.9 (307.0)	36.42	1098.3 (999.3)	24.18	1411.3 (1333.8)	23.88	1382.6 (1295.7)

Table A7: Complete accuracy performance and inference efficiency (unit: words) of DeepSeek-R1-Distill model family on vertical analysis (temporal reasoning). The numbers in parentheses indicate the thought lengths.

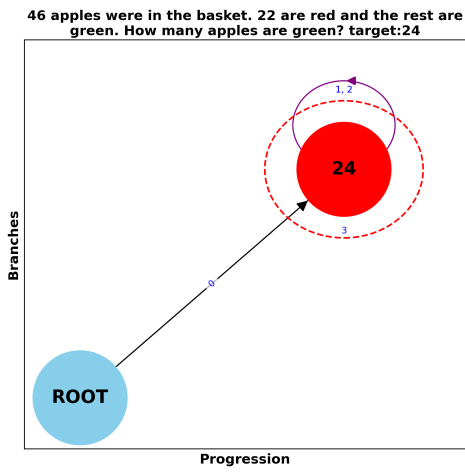


Figure A6: Individual thought progression graph of Qwen3-235B-A22B model on a sampled grade-1 arithmetic query.

What is the time 2316 day after Nov 18, 1185 target:March 22, 1192

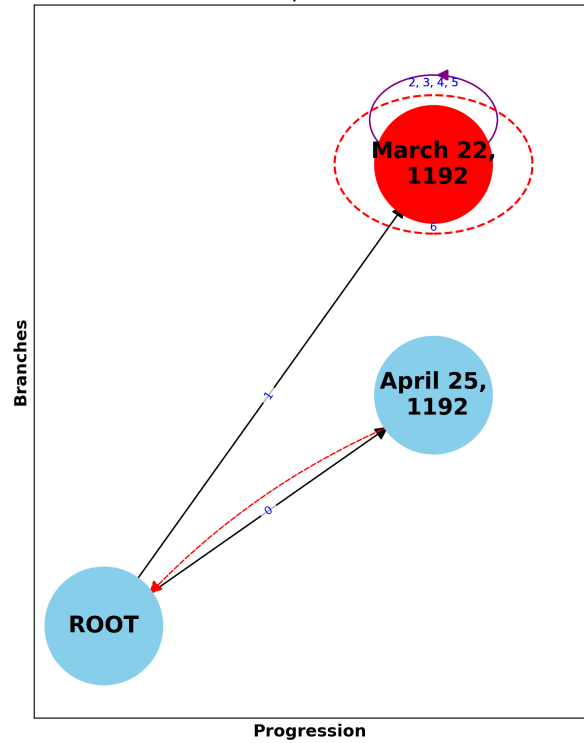


Figure A8: Individual thought progression graph of Qwen3-32B model on a sampled date arithmetic (temporal-L3) query.

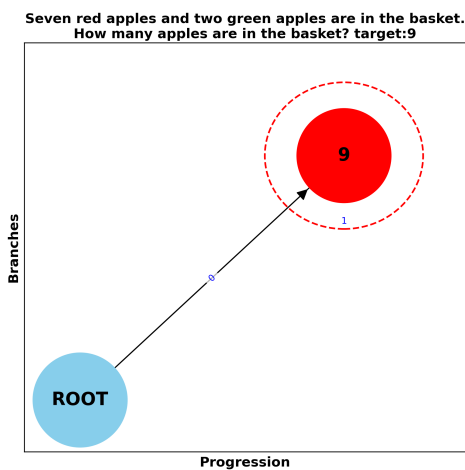


Figure A7: Individual thought progression graph of R1-Distill-Llama-70B model on a sampled grade-1 arithmetic query.

Seven red apples and two green apples are in the basket. How many apples are in the basket? target:9

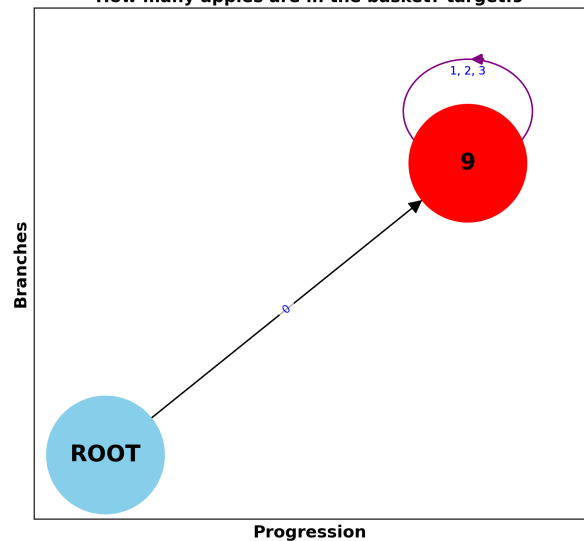


Figure A9: Individual thought progression graph of Qwen3-32B model on a sampled grade-1 arithmetic query.

Jim has a 20 pack of gum. He chews 1 piece of gum for every 2 hours he's at school over a school day that lasts 8 hours. He chews 1 piece on the way home from school and 1 stick after dinner. He also gives half the gum he has remaining to his sister when she asks for some right before bed. How many pieces of gum does Jim have left at the end of the day? target:7

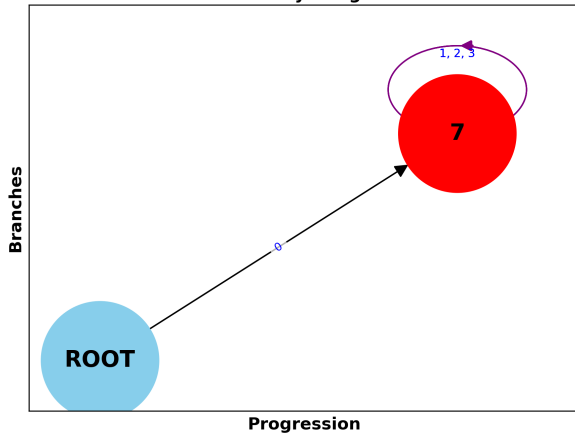


Figure A10: Individual thought progression graph of Qwen3-32B model on a sampled GSM8k query.

A marketing company pays its employees on a commission-based salary system. If you sell goods worth 1000, you earn a 30% commission. Sales over 1000 get you an additional 10% commission. Calculate the amount of money Antonella earned if she sold goods worth 2500. target:450

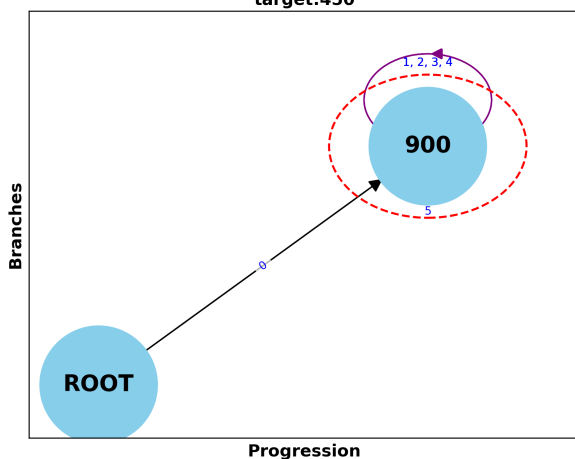


Figure A11: Individual thought progression graph of Qwen3-32B model on a sampled GSM8k query.

What is the time 2316 day after Nov 18, 1185 target:March 22, 1192

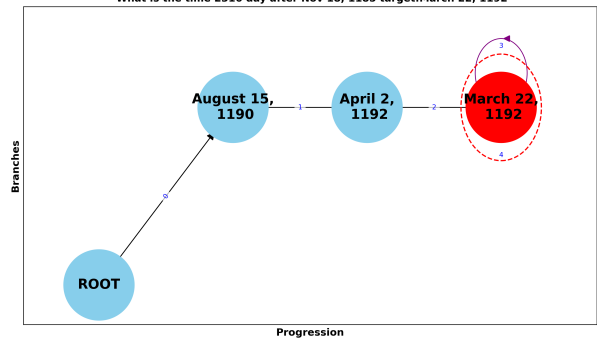


Figure A12: Individual thought progression graph of R1-Distill-Llama-70B model on a sampled date arithmetic (temporal-L3) query.

Doctor Jones is scheduling his time for Monday. He is spending nine hours at the clinic where he works that day. He has to do rounds to check on inpatients staying at the clinic, which takes twenty minutes per inpatient, and he has ten appointments, which take thirty minutes each. How many hours will Doctor Jones have left to update his records if he has 9 inpatients at the clinic? target:1

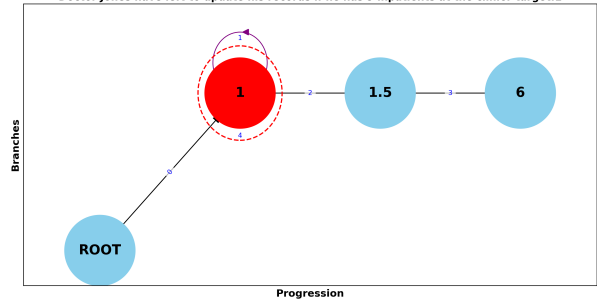


Figure A13: Individual thought progression graph of R1-Distill-Llama-70B model on a sampled GSM8k query.

Jim has a 20 pack of gum. He chews 1 piece of gum for every 2 hours he's at school over a school day that lasts 8 hours. He chews 1 piece on the way home from school and 1 stick after dinner. He also gives half the gum he has remaining to his sister when she asks for some right before bed. How many pieces of gum does Jim have left at the end of the day? target:7

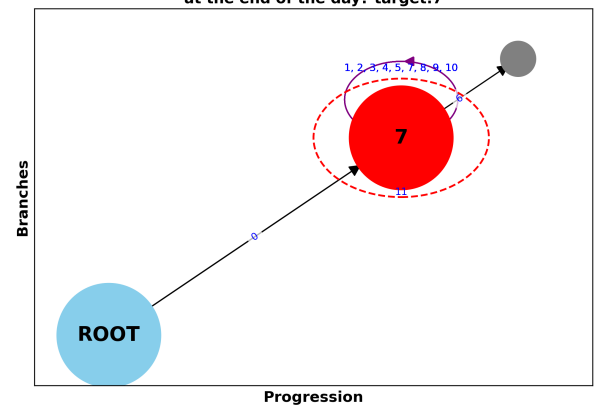


Figure A14: Individual thought progression graph of Qwen3-235B-A22B model on a sampled GSM8k query.

A marketing company pays its employees on a commission-based salary system. If you sell goods worth 1000, you earn a 30% commission. Sales over 1000 get you an additional 10% commission. Calculate the amount of money Antonella earned if she sold goods worth 2500. target:450

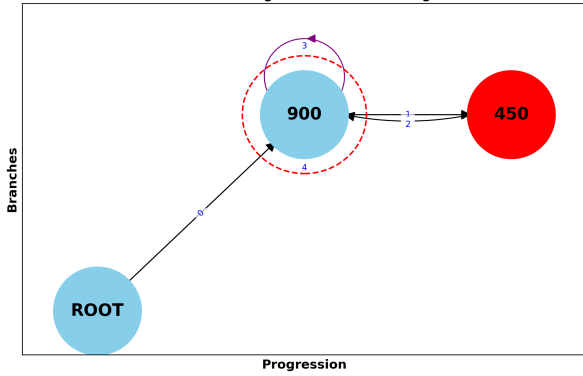


Figure A15: Individual thought progression graph of Qwen3-235B-A22B model on a sampled GSM8k query.

Doctor Jones is scheduling his time for Monday. He is spending nine hours at the clinic where he works that day. He has to do rounds to check on inpatients staying at the clinic, which takes twenty minutes per inpatient, and he has ten appointments, which take thirty minutes each. How many hours will Doctor Jones have left to update his records if he has 9 inpatients at the clinic? target:1

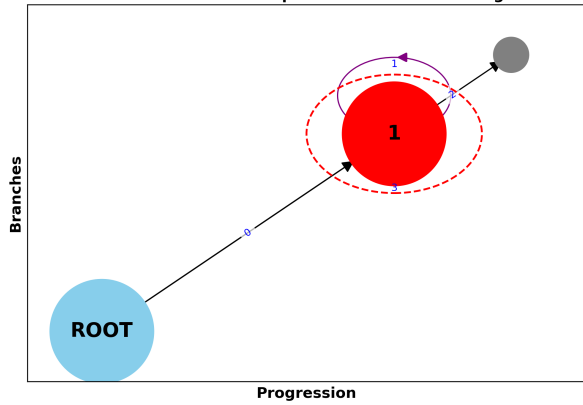


Figure A16: Individual thought progression graph of Qwen3-235B-A22B model on a sampled GSM8k query. This is the graphical representation of Figure A34 and A35.

Induced 2-answer Thought Progression Pattern of Qwen3-235B-A22B on Temporal-L3 (Occurrence: 25.9%)

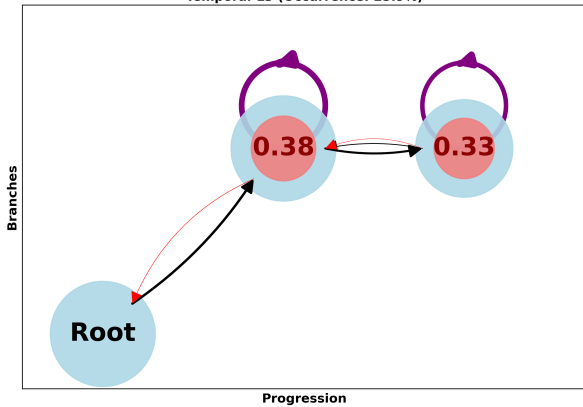


Figure A17: The **Explorer** thought progression pattern (2 distinct answer case). This is a trivial case.

Induced 3-answer Thought Progression Pattern of Qwen3-235B-A22B on Temporal-L3 (Occurrence: 26.2%)

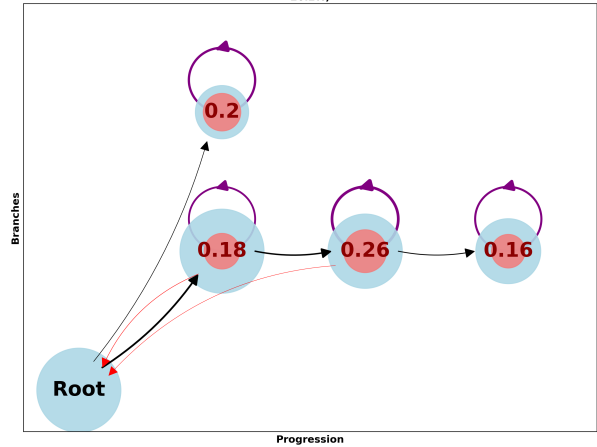


Figure A18: The **Explorer** thought progression pattern (3 distinct answer case).

Induced 4-answer Thought Progression Pattern of Qwen3-235B-A22B on Temporal-L3 (Occurrence: 13.6%)

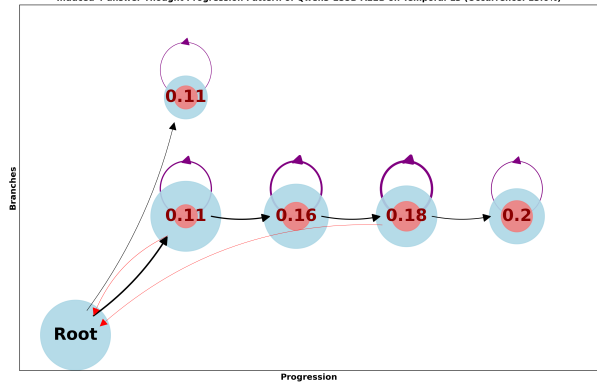


Figure A19: The **Explorer** thought progression pattern (4 distinct answer case).

Induced 6-answer Thought Progression Pattern of Qwen3-235B-A22B on Temporal-L3 (Occurrence: 2.4%)

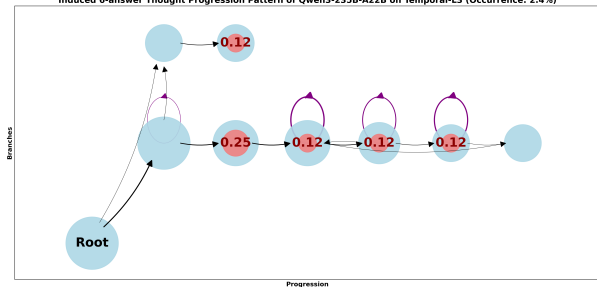


Figure A20: The **Explorer** thought progression pattern (6 distinct answer case).

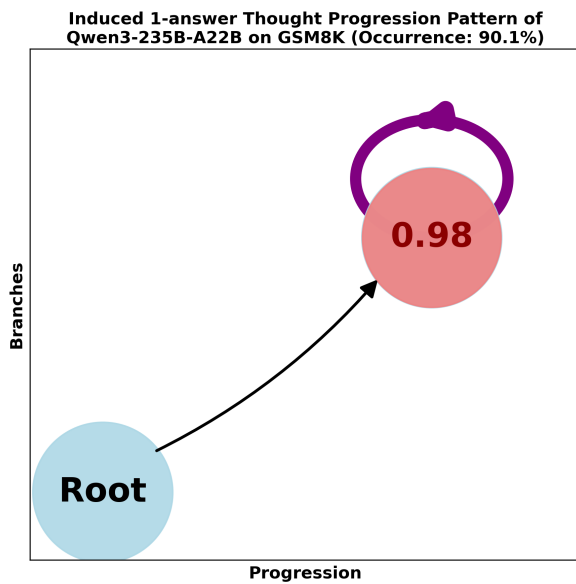


Figure A21: The **Explorer** thought progression pattern (1 distinct answer case). This is a trivial case.

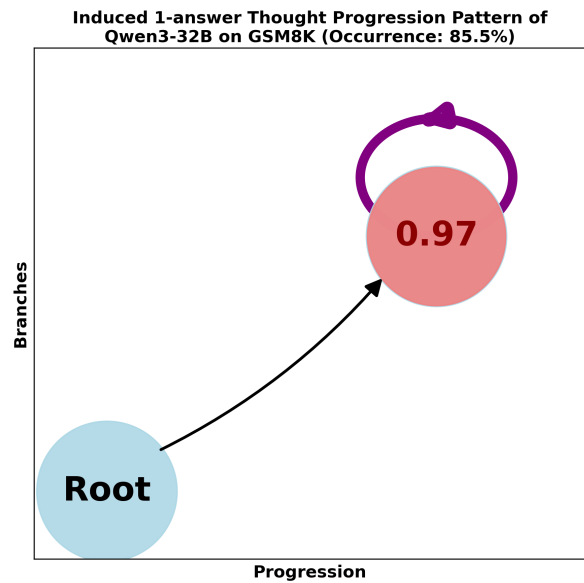


Figure A23: The **Late Landing** thought progression pattern (1 distinct answer case). This is a trivial case.

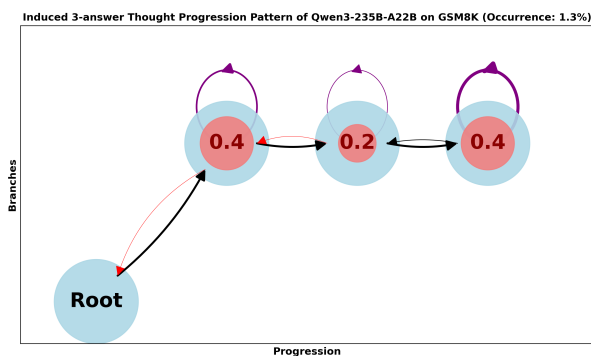


Figure A22: The **Explorer** thought progression pattern (3 distinct answer case).

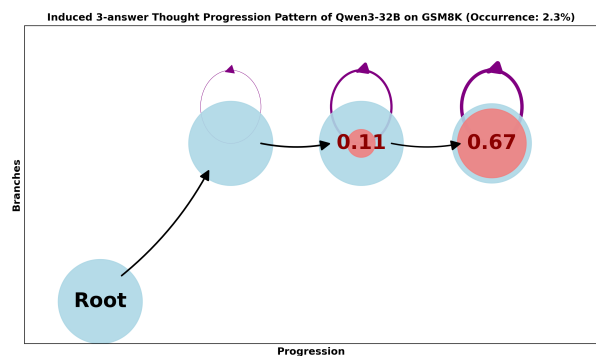


Figure A24: The **Late Landing** thought progression pattern (3 distinct answer case).

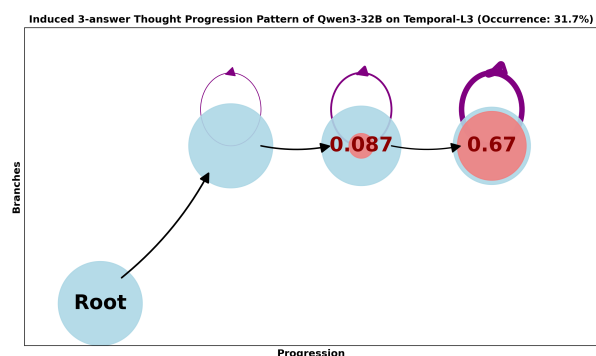


Figure A25: The **Late Landing** thought progression pattern (3 distinct answer case).

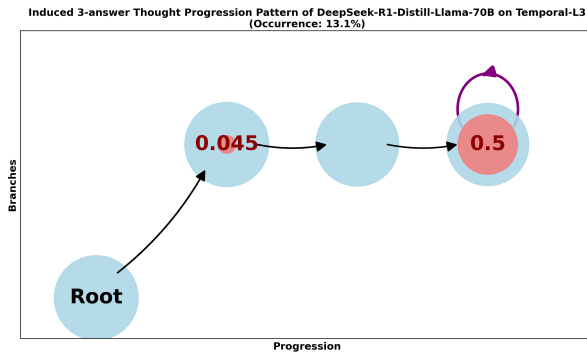


Figure A26: The **Late Landing** thought progression pattern (3 distinct answer case).

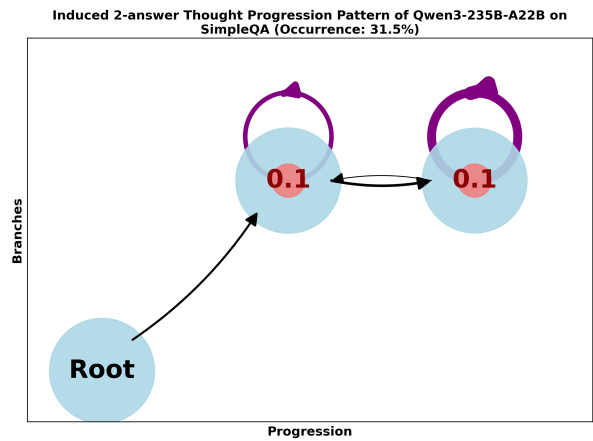


Figure A29: The **Explorer** thought progression pattern (2 distinct answer case) on the SimpleQA dataset which requires minimal reasoning efforts.

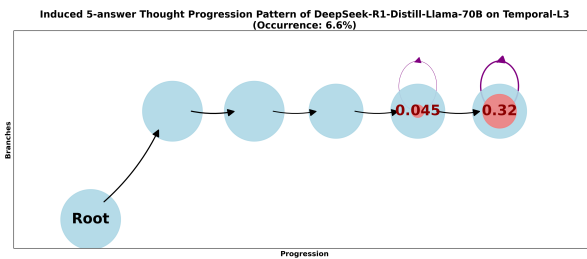


Figure A27: The **Late Landing** thought progression pattern (5 distinct answer case).

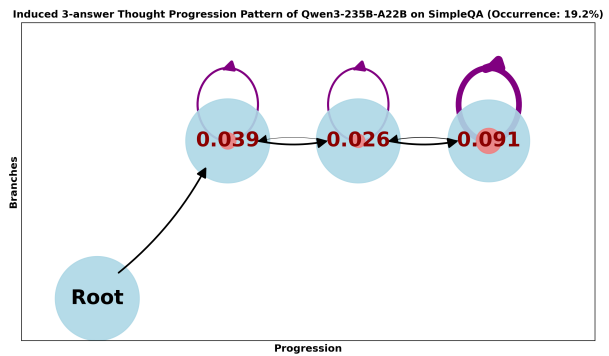


Figure A30: The **Explorer** thought progression pattern (3 distinct answer case) on the SimpleQA dataset which requires minimal reasoning efforts.

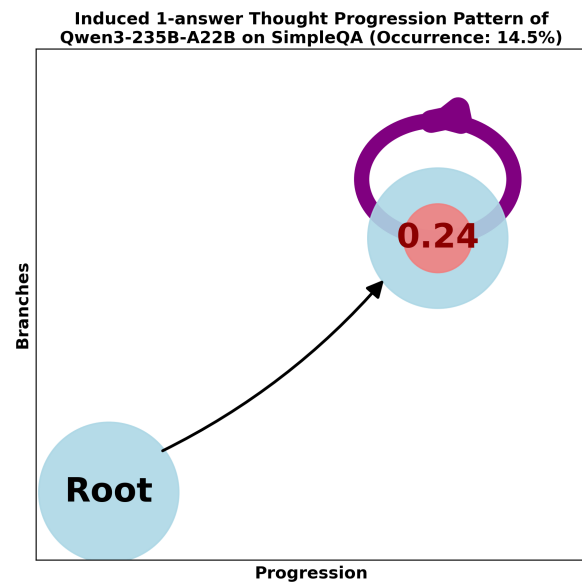


Figure A28: The **Explorer** thought progression pattern (1 distinct answer case) on the SimpleQA dataset which requires minimal reasoning efforts.

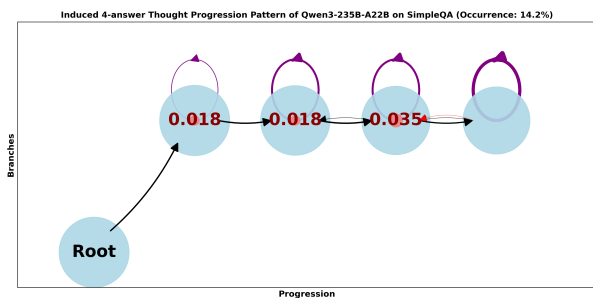


Figure A31: The **Explorer** thought progression pattern (4 distinct answer case) on the SimpleQA dataset which requires minimal reasoning efforts.

Instructions for Sub-Thought Decomposition and Label Inference (part 1)

Your goal is to meticulously analyze a machine's thought process, breaking it down into individual sub-thoughts to understand its reasoning flow. This detailed breakdown helps in evaluating and improving AI transparency and logic.

Here's the input format you'll receive:

Context: {context}

Question: {question}

Gold target: {target}

Thought process: {thought_process}

Your Task: Sub-Thought Decomposition

1. **Decompose the Thought Process:** Break down the entire machine-generated Thought process into individual sub-thoughts. Each sub-thought must be:

- **Self-contained:** It must make sense on its own, without heavy reliance on other sub-thoughts.
- **Complete:** A sub-thought must represent a **complete** sequence of logical steps, generally consisting of intent, action/calculation, and conclusion in that order.
- **Answer-bearing:** A proposed answer must be **clearly** visible in each sub-thought, which should be in the same format as the Gold target.

Determine the boundary of each sub-thought by locating **pivoting phrases**. These are linguistic cues that indicate a *shift* in the thought process. More importantly, pivoting phrases is *always* at the **beginning** of a sub-thought.

Common pivoting phrases include but are not limited to:

- Wait, Alternatively, Another angle, Another approach, But wait, However, Hold on, Let me double-check, On the other hand

Crucial Rule 1: Avoid Isolating Pivoting Phrases

A pivoting phrase **cannot be a sub-thought by itself**. It is a signal of a new thought, not the thought itself. The sub-thought **must include the action or reasoning** that the pivot introduces.

Crucial Rule 2: Avoid Fragmented Intent

A statement of intent (e.g., "I will now calculate...") and its immediate execution must be treated as a single, indivisible sub-thought. The intent and the action are two parts of the same logical step and should not be split.

Crucial Rule 3: Avoid Fragmented Conclusions

A sub-thought must represent a complete logical step. A concluding phrase (e.g., "So, the answer is...", "Therefore, the result is...") should **NOT** be separated from the immediately preceding calculation, reasoning, or verification that produced it. The reasoning and its conclusion are part of the **same** sub-thought.

2. **Analyze Each Sub-Thought:** For every sub-thought you identify, provide:

- A **short description** summarizing how the sub-thought approaches the question.
- Its **discourse attribute**, classifying its role in the overall thought process.

Figure A32: Instructions for Sub-Thought Decomposition and Label Inference (part 1). Part 2 is shown in Figure A33.

Instructions for Sub-Thought Decomposition and Label Inference (part 2)

Discourse Attributes Defined:

- **initial:** The very first attempt of solving the question. The initial sub-thought might include an outline of the overall plan. There is only one initial sub-thought.
- **verification:** Verify and confirms the correctness of the *immediately preceding* sub-thought. No changes or corrections are made.
- **correction:** Modifies or rectifies the *immediately preceding* sub-thought.
- **backtrack:** Abandons the current line of reasoning and reverts to an earlier (*NOT* immediately preceding) sub-thought to pursue a different path from that point.
- **sidetrack:** Explores new, supplementary information that isn't directly related to a different approach to the main question.
- **branching out:** Explores alternative ways/methods to answer the user's question or different aspects of it, as opposed to simply adding new information (sidetrack).
- **final:** The final sub-thought that provides the final proposed answer.

Note: For attributes that inherently relate to a specific previous sub-thought, use the following format: '[Attribute Type]->[Target Index]'.

- This *ONLY* applies to **verification**, **correction**, and **backtrack**
- The [Target Index] is the 1-indexed position of the sub-thought being referred to.
- For verification and correction, the target will *ALWAYS* be the immediately preceding sub-thought (e.g., a correction at index **4** would be formatted as 'correction->3').
- For backtrack, the target can be any earlier, *NON-PRECEDING* sub-thought (e.g., a backtrack at index **5** that reverts to the thought at index **1** would be formatted as 'backtrack->1'). If there is a "start over" action, mark it as 'backtrack->ROOT'.

Expected Output Format:

Your analysis output must be a JSON array structured as follows:

```
[
  {
    "sub_thought_index": "The 0-indexed position of the sub-thought",
    "sub_thought_description": "A concise summary of what this
sub-thought is about",
    "sub_thought_answer": "The machine-generated answer extracted
from this *CURRENT* sub-thought. It should be in the same format
as the ``Gold target''",
    "discourse_attribute": "The classified role of this sub-thought
in the entire thought process. Each sub-thought can only have
one discourse attribute. Use your best judgement to
determine the most appropriate attribute.",
    "sub_thought_content": "The exact text segment from the original
thought process that constitutes this sub-thought"
  }
]
```

Figure A33: Instructions for Sub-Thought Decomposition and Label Inference (part 2).

Sample Parsed Output by TRACE (Sub-thoughts 0-1)

Sub-thought 0

Content: Okay, let's see. Doctor Jones has a total of nine hours at the clinic on Monday. He needs to spend time on rounds for inpatients and ten appointments. Then he has to update his records. The question is asking how many hours he has left for updating records.

First, I need to calculate the total time he spends on rounds and appointments. Then subtract that from the nine hours to find the remaining time for records.

Let me start with the rounds. He has 9 inpatients, and each takes 20 minutes. So 9 times 20 minutes. Let me compute that: $9 * 20 = 180$ minutes.

Then the appointments. He has ten appointments, each 30 minutes. So $10 * 30 = 300$ minutes.

Total time spent on both rounds and appointments is $180 + 300$ minutes. That's 480 minutes. Now, convert that to hours because the total time he has is in hours. Since 60 minutes = 1 hour, divide 480 by 60. $480 \div 60 = 8$ hours.

Wait, so he spends 8 hours on those two tasks. He has nine hours total. So $9 - 8 = 1$ hour left for updating records. That seems straightforward.

Discourse Attribute: initial

Sub-thought Answer:

Generated Answer:

Sub-thought 1

Content: Wait, but let me double-check. Maybe I made a mistake in the calculation. Let me go through again.

Inpatient rounds: 9 patients * 20 minutes each. $9 * 20$ is 180. Correct. Appointments: $10 * 30 = 300$. Correct. Total $180 + 300 = 480$. 480 divided by 60 is 8. So yes, 8 hours. Subtract from 9 hours, 1 hour left. So the answer should be 1 hour.

Discourse Attribute: verification

Sub-thought Answer (Correctness):

Answer in Response (Helpfulness):

Figure A34: Sample parsed output (sub-thoughts 0-1) by TRACE on the Qwen3-235B-A22B-generated response to a GSM8K prompt. The prompt is "Doctor Jones is scheduling his time for Monday. He is spending nine hours at the clinic where he works that day. He has to do rounds to check on inpatients staying at the clinic, which takes twenty minutes per inpatient, and he has ten appointments, which take thirty minutes each. How many hours will Doctor Jones have left to update his records if he has 9 inpatients at the clinic?". The ground-truth answer is 1. The constructed thought progression graph for this thought process is shown in Figure A16.

Sample Parsed Output by TRACE (Sub-thought 2-3)

Sub-thought 2

Content: Hmm. But wait, the problem says "how many hours", so the answer is 1 hour. But maybe I should check if there's any hidden details. Like, maybe the appointments and rounds are scheduled during the nine hours, but does that include any breaks or not? The problem doesn't mention any breaks, so we can assume that the nine hours are entirely his time to be scheduled, and he has to allocate it to rounds, appointments, and updating records.

Discourse Attribute: sidetrack

Sub-thought Answer (Correctness):

Answer in Response (Helpfulness):

Sub-thought 3

Content: Therefore, the answer should be 1 hour. Let me just make sure that the conversion from minutes to hours is correct. 480 minutes divided by 60 is 8 hours. Yep. $9 - 8 = 1$. I think that's right.

Discourse Attribute: final

Sub-thought Answer (Correctness):

Answer in Response (Helpfulness):

Figure A35: Sample parsed output (sub-thoughts 2-3) by TRACE on the Qwen3-235B-A22B-generated response to a GSM8K prompt. Sub-thought 2 is marked as *sidetrack* since the model is rambling and exploring supplementary information which does not directly contribute to solve the input problem. The graphical representation of *sidetrack* is a tiny, gray bubble as shown in Figure A16.