

# Can LLMs Reason Abstractly Over Math Word Problems Without CoT? Disentangling Abstract Formulation From Arithmetic Computation

Ziling Cheng<sup>1,2†</sup> Meng Cao<sup>1,2</sup>

Leila Pishdad<sup>3</sup> Yanshuai Cao<sup>3</sup> Jackie Chi Kit Cheung<sup>1,2,4</sup>

<sup>1</sup>Mila – Quebec AI Institute <sup>2</sup>McGill University <sup>3</sup>RBC Borealis

<sup>4</sup>Canada CIFAR AI Chair

{ziling.cheng, meng.cao}@mail.mcgill.ca

{leila.pishdad, yanshuai.cao}@borealisai.com, cheungja@mila.quebec

## Abstract

Final-answer-based metrics are commonly used for evaluating large language models (LLMs) on math word problems, often taken as proxies for reasoning ability. However, such metrics conflate two distinct sub-skills: **abstract formulation** (capturing mathematical relationships using expressions) and **arithmetic computation** (executing the calculations). Through a disentangled evaluation on GSM8K and SVAMP, we find that the final-answer accuracy of Llama-3 and Qwen2.5 (1B-32B) without CoT is overwhelmingly bottlenecked by the arithmetic computation step and not by the abstract formulation step. Contrary to the common belief, we show that CoT primarily aids in computation, with limited impact on abstract formulation. Mechanistically, we show that these two skills are composed conjunctively even in a single forward pass without any reasoning steps via an **abstract-then-compute** mechanism: models first capture problem abstractions, then handle computation. Causal patching confirms these abstractions are present, transferable, composable, and precede computation. These behavioural and mechanistic findings highlight the need for disentangled evaluation to accurately assess LLM reasoning and to guide future improvements.<sup>1</sup>

## 1 Introduction

Large language models (LLMs) have demonstrated impressive progress on various math problem datasets (Cobbe et al., 2021; Hendrycks et al., 2021b; Patel et al., 2021), often leveraging Chain-of-Thought (CoT) prompting (Wei et al., 2022). Despite the availability of step-by-step reasoning chains, standard evaluation predominantly relies on final-answer accuracy (comparing the model’s final numerical output against a gold answer), which

reduces model performance to a single metric (Liu et al., 2024; Opedal et al., 2024). This reduction limits the possible insights when diagnosing LLMs’ reasoning abilities, especially in zero-shot scenarios without CoT. When an LLM fails to produce the correct answer, is it due to “reasoning deficits”, or could it be a calculation error?

To investigate this, we propose a disentangled evaluation framework that separately measures two core skills of mathematical problem-solving (See Figure 1): (1) **abstract formulation** (hereafter, abstraction) — the ability to identify relevant quantities and translate the natural language problem into its underlying mathematical relationships (e.g.,  $36 + 47$  or  $x + y$  in Figure 1); and (2) **arithmetic computation** (hereafter, computation) — the capacity to calculate the final answer from that expression (e.g., evaluate  $36 + 47$  to 83).

Using this disentangled evaluation on GSM8K (Cobbe et al., 2021) and SVAMP (Patel et al., 2021) with Llama-3 and Qwen-2.5 models (1B-32B), we find that even without CoT: (i) models surprisingly perform better at abstraction than computation, despite the former’s perceived conceptual complexity. (ii) if deriving the final answer in math word problems depends on these two skills conjunctively, final-answer accuracy alone may give a misleading picture of models’ reasoning abilities in math word problems. Moreover, we show that CoT primarily improves computation, with limited gains in abstraction, further demonstrating the value of disentangled evaluation.

While these behavioural findings suggest that models can formulate abstractions without explicit CoT when separately prompted, it remains unclear whether abstraction and computation are composed conjunctively when deriving the final answer during single-pass inference. To explore this, we move beyond outcome-based evaluation, and conduct mechanistic interpretability analyses. Using logit attribution and activation patching, we identify a

<sup>†</sup>Work done during a Mitacs internship at Borealis AI.

<sup>1</sup>Code available at: <https://github.com/ziling-cheng/Disentangle-Math-Reasoning>.

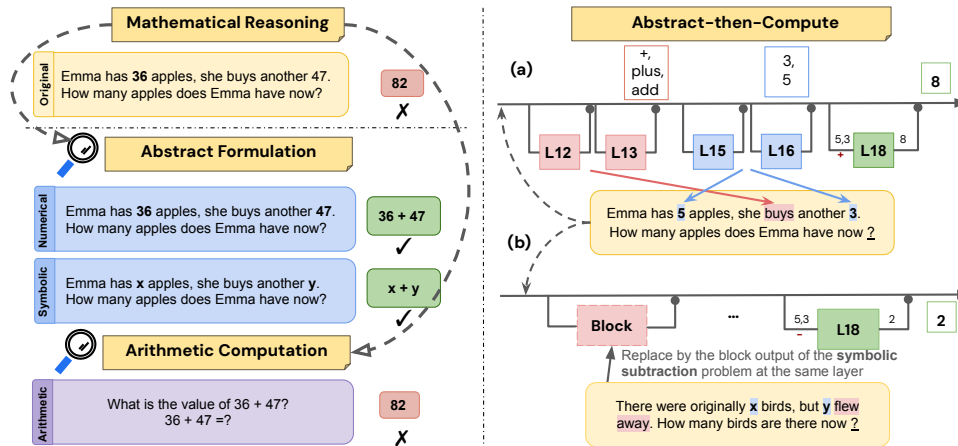


Figure 1: **Left (Disentangled evaluation framework):** Final-answer accuracy obscures reasoning ability due to conflating abstract formulation and arithmetic computation. **Right (Abstract-then-Compute Mechanism in Llama-3 8B):** (a) Residual stream at the last token position shows that models first capture problem abstraction (L13-14), followed by computation (L18). (b) Same as (a), but one critical layer output is patched with a different symbolic abstraction (e.g.,  $x - y$ ), causally changing the computation from  $5 + 3 = 8$  to  $5 - 3 = 2$ .

consistent and sequential **abstract-then-compute** mechanism (see Figure 1a), where models first capture problem abstraction (L13-14), followed by arithmetic computation (L18). Moreover, cross-prompt patching provides evidence that models do form abstractions internally independent of the surface form (numerical or symbolic, see Figure 1b): when these symbolic abstractions (e.g.  $x - y$ ) are transferred into a different problem, they are utilized and composed with the subsequent computation stages, altering the final answer.

**Contributions:** (i) Through disentangled evaluation, we show that without CoT, models exhibit stronger reasoning ability than final-answer accuracy suggests, and that CoT primarily aids calculation. (ii) Using mechanistic interpretability, we uncover an abstract-then-compute mechanism in a single-pass generation, where abstractions are transferrable across problem variants. Collectively, our findings suggest an alternative narrative: poor final-answer accuracy without CoT (Wei et al., 2022; Sprague et al., 2025), or performance declines on problem variants (Zhang et al., 2024a; Shi et al., 2023; Mirzadeh et al., 2025), can stem from arithmetic errors rather than reasoning deficits.

## 2 Related Work

**Mathematical Reasoning Evaluation** Existing math problem-solving benchmarks spans elementary word problems (Cobbe et al., 2021; Patel et al., 2021; Amini et al., 2019; Miao et al., 2020; Ling et al., 2017; Koncel-Kedziorski et al., 2016; Shi

et al., 2015) to higher levels (Hendrycks et al., 2021b,a; Zhong et al., 2024; Zhang et al., 2023; He et al., 2024). Early datasets paired expressions with answers, but evaluation largely focused on final-answer-based metrics (Patel et al., 2021; Shi et al., 2015). With the rise of LLMs and CoT prompting (Wei et al., 2022), rationale-based formats became common (Hendrycks et al., 2021b; Cobbe et al., 2021), with some approaches rely on program generation and execution (Mishra et al., 2022; Gao et al., 2023; Chen et al., 2023), yet standard evaluations still predominantly use final-answer metrics. In contrast, we move beyond this final-answer-centric paradigm, by decomposing problem-solving into abstract formulation and arithmetic computation, inspired by the cognitive theories (Opedal et al., 2024).

**Memorization vs. Generalization** Variants of math word problems with perturbations were introduced to test generalization beyond memorization (Zhang et al., 2024a; Ye et al., 2025; Gao et al., 2023; Shi et al., 2023; Li et al., 2024; Mirzadeh et al., 2025). While performance drops are often interpreted as reasoning failures, our results suggest they may instead stem mainly from arithmetic errors, pointing to a different improvement strategy.

**Mechanistic Interpretability** Mechanistic interpretability methods, such as logit attribution (nostalgebraist, 2020; Belrose et al., 2023) and causal patching (Goldowsky-Dill et al., 2023; Wang et al., 2023; Meng et al., 2022; Zhang and Nanda, 2023; Merullo et al., 2024; Cheng et al., 2025), have

been used to trace model computations. Most prior mechanistic interpretability work in math reasoning has focused on synthetic arithmetic calculation tasks (e.g.,  $226 - 68 = ?$ ), where models are prompted to perform explicit calculations (Nikankin et al., 2025; Zhang et al., 2024b; Yu and Ananiadou, 2024; Zhou et al., 2024). For example, Nikankin et al. (2025) revealed an arithmetic circuit with activation patching and classifier probing: attention heads copy operands and operators, and MLPs compute the answer using heuristics. Zhou et al. (2024) identified that LLMs add numbers using Fourier features with logit lens. Yu and Ananiadou (2024) proposed the Comparative Neuron Analysis (CNA) method and identified four distinct stages in the internal logic chain (feature enhancing, feature transferring, feature predicting, and prediction enhancing) for calculating numbers. In contrast, we focused on a more complex setup – math word problems – which requires more processing from natural language. While (Ye et al., 2025) applied probing classifiers to identify surface-level answer attributes (e.g., parameters computed and to compute), our work applies existing mechanistic interpretability tools to understand if the two stages (abstraction and computation) are implemented by the models internally in a sequential manner. We uncover an abstract-then-compute mechanism in math word problems, and we further show that these abstract representations are transferable across different problems and surface forms — providing new insight into the internal reasoning in LLMs.

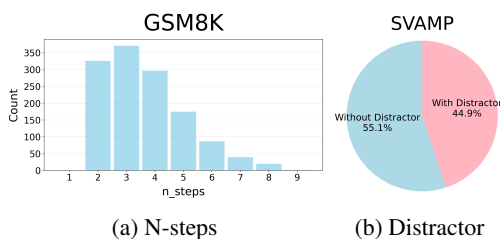


Figure 2: Distribution of problem characteristics by number of reasoning steps (GSM8K) and presence of distractors (SVAMP).

### 3 Dataset and Experimental Design

**Task and Dataset** We study math word problems using GSM-8K (Cobbe et al., 2021) and SVAMP (Patel et al., 2021). GSM-8K spans 2–8 steps without distractors (See Figure 2 for statistics), while SVAMP involves single-step reasoning with dis-

tractor variants. To evaluate abstract formulation, we create symbolic variants: SVAMP expressions are templated into variable-based forms; GSM-8K symbolic versions from the test set are generated using gpt-4o-mini (OpenAI, 2024) via a two-stage *generate-then-validate* to ensure the correctness. See Appendix A.1 and Table 4 for details and examples. For interpretability, we generate 3,600 simple 1–2 step<sup>2</sup> word problems involving basic operations (+, −, ×, ÷) from 1,200 diverse LLM-generated templates, covering varied scenarios, verb choices, entities, names and sentence structures. See Appendix B.1 and Table 7 for details and examples.

**Models** We evaluate instruction-tuned Llama-3 (1B, 3B, 8B) (Grattafiori et al., 2024) and Qwen 2.5 (Yang et al., 2024) (3B, 7B, 14B, 32B) models. Mechanistic interpretability analyses focus on Llama-3 8B, Qwen 2.5 7B, and Qwen 2.5 14B.

**Evaluation** All experiments use greedy decoding and FP16 precision on RTX 8000/A100L GPUs. Numeric answers are evaluated via normalized Exact Match. Symbolic expressions are evaluated using gpt-4o-mini (94% agreement with humans on 120 samples, prompt and details in Appendix A.2) and with sympy for numeric expressions. We report standard accuracy. CoT generations are capped at 512 tokens. See Appendix A.2 for details.

### 4 Disentangled Evaluation

We first introduce the disentangled evaluation framework, then present results without CoT in Sec. 4.1, followed by an analysis of CoT’s impact in Sec. 4.2.

**Framework** Suppose a task  $T$  can be decomposed into a set of sub-skills  $\{s_1, s_2, \dots, s_n\}$ , such that solving  $T$  requires executing these skills conjunctively (i.e.,  $T = s_1 \cap s_2 \cap \dots \cap s_n$ ). Disentangled evaluation aims to assess each sub-skill  $s_i$  independently via a corresponding subtask  $t_i$ , designed to isolate and test that specific skill. Let  $\text{Eval}(T)$  denote the evaluation metric on the full task, and  $\text{Eval}(t_i)$  the metric for subtask  $t_i$ . Measuring  $\text{Eval}(t_1), \dots, \text{Eval}(t_n)$  enables finer-grained attribution of performance, identifying failure of specific skills. In math word problems, let  $(Q, E, A)$  be the question, expression and answer

<sup>2</sup>We focus on 1–2 step problems, as models often fail simple word problems involving multi-step computations in a single forward pass.

Setting	Skills Tested	Question Form and Example	Answer Form and Example
<b>Original</b>	Abstraction +Computation	<u>Numerical</u> : Weng earns \$12 for every hour she works. If she worked for 50 minutes, how much did she earn?	<u>Number</u> : 10
<b>Arithmetic Computation</b>	Computation	<u>Numerical</u> : What is the value of $12 \times (\frac{50}{60})$ ?	<u>Number</u> : 10
<b>Numerical Abstraction</b>	Abstraction	<u>Numerical</u> : Weng earns \$12 for every hour she works. If she worked for 50 minutes, how much did she earn?	<u>Expression</u> : $12 \times (\frac{50}{60})$
<b>Symbolic Abstraction</b>	Abstraction	<u>Symbolic</u> : Weng earns \$x for every hour she works. If she worked for y minutes, how much did she earn?	<u>Expression</u> : $x \times (\frac{y}{60})$

Table 1: Disentangled evaluation in math word problems with tested skills, varying by question and answer forms. Instructions in Appendix Table 5.

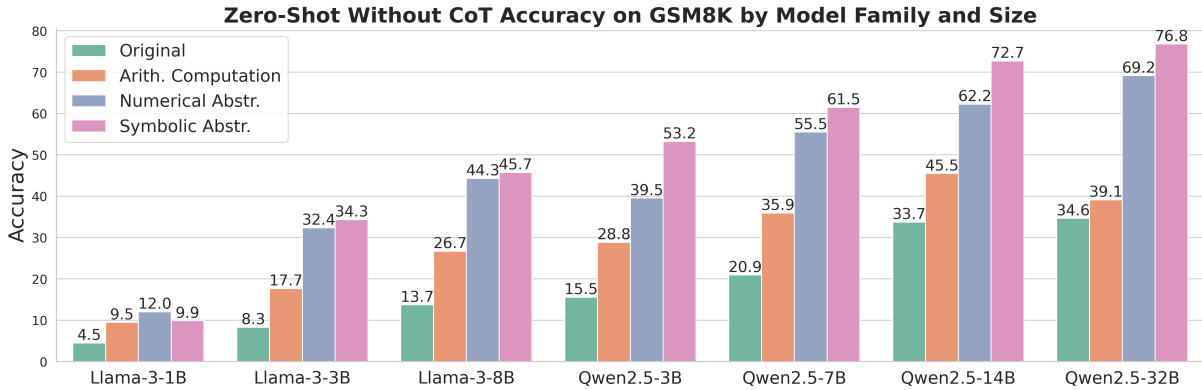


Figure 3: Model zero-shot **without CoT** performance on GSM8K. (i) Models exhibit much better abstraction performance (*Symbolic* and *Numerical*) than in actually computing the expressions (*Arithmetic Computation*). (ii) Final-answer accuracy in the *Original* setting may provide a misleading picture of models’ reasoning ability, possibly due to arithmetic limitations.

triplets, we decompose mathematical problem-solving into **abstract formulation** (translating  $Q$  to mathematical relationships  $E$ ) and **arithmetic computation** (executing the calculation from  $E$  to produce  $A$ ). Besides the standard *Original* setting (requiring both abstraction and computation), we design three targeted subtasks: *Symbolic Abstraction*, which assess abstraction using symbolic variables; *Numerical Abstraction*, evaluating abstraction with concrete numbers but without computation; and *Arithmetic Computation*, which directly tests execution of fully specified expressions from  $Q$ . See Table 1 and Appendix A.2 for details.

#### 4.1 Understanding Model Failures: Reasoning or Arithmetic Error?

We first apply disentangled evaluation **zero-shot without CoT** across multiple model sizes of Llama-3 and Qwen2.5 families. As shown in Figure 3, the error rates are consistently lower for abstract formulation (both *Numerical* and *Symbolic Abstraction*) compared to arithmetic computation. This suggests that if final-answer accuracy in the *Original* set-

ting depends on both competencies conjunctively, poor performance observed in the *Original* setting could stem from arithmetic computation failures, rather than reasoning deficits. Consequently, this indicates that final-answer accuracy alone from the *Original* setting may substantially mislead a model’s underlying reasoning ability. See additional results in Appendix A.3. To assess the reliability and external validity of the symbolic abstraction evaluation, we perform ablations over symbol order and symbol choice in Appendix A.5.

$\Delta$ Accuracy	8B	7B	14B	32B	Avg.
Original	64.8	68.5	58.4	59.7	62.8
Arith. Comp.	64.8	60.5	51.2	58.2	58.7
Numerical Abstr.	15.8	21.6	21.6	11.6	17.6
Symbolic Abstr.	11.0	13.2	1.1	1.3	6.7

Table 2: Accuracy difference (%) with and without CoT. Results are shown for Llama 3 (8B) and Qwen2.5 models (7B, 14B, 32B).

## 4.2 Disentangling CoT Gains

We now apply disentangled evaluation with CoT to disentangle CoT gains (Table 2). We show that CoT yields the largest gains in computation (e.g., +62.8%), confirming its effectiveness in multi-step arithmetic. In contrast, abstraction shows limited improvement (e.g., +6.7% for *Symbolic abstraction* and +17.6% for *Numerical abstraction*), even with extended generation budgets (512 tokens), suggesting CoT is less helpful for abstraction. Gains in the *Original* setting (e.g., +62.8%) likely reflect a mix of benefits from both components and possible data leakage. See additional results in Appendix A.4.

**Summary:** These findings challenge the view that poor final-answer accuracy in math reasoning benchmarks always implies ‘poor reasoning’. Instead, our disentangled design reveals that many models do possess a level of abstract formulation capabilities, which are often obscured in standard evaluations due to their limited arithmetic competence. Crucially, while abstraction variants indicate far higher performance than the *Original* setting, models are still not perfect — performance in *Symbolic Abstraction* remains far from 100% (45.7% for Llama-8B, 76.8% for Qwen-32B), but the gap is significantly narrower than previously assumed, calling for more precise definition and evaluation of reasoning.

## 5 Inside the Model: Probing Abstraction and Computation

To investigate whether abstraction and computation are composed conjunctively when producing a final numerical answer in a single forward pass, we move beyond outcome-based evaluation and apply mechanistic interpretability. We hypothesize an **abstract-then-compute** process: first inferring the abstraction (e.g., ‘+’ from “buys”), then performing the computation (e.g.,  $5 + 3$ ). Section 5.1 identifies key layers for each stage; Section 5.2 validate these layers and tests abstraction transferability across forms (symbolic/concrete) and logic.

### 5.1 Uncovering the Abstract-Then-Compute Mechanism in One Forward Pass

#### 5.1.1 Methods

We use logit attribution (nostalgebraist, 2020; Belrose et al., 2023) and activation patching (Ghandeharioun et al., 2024; Zhang and Nanda, 2023; Meng et al., 2022) to probe whether abstraction and computation occur during single-step genera-

tion. As summarized in Figure 4, we seek evidence of abstraction and computation.

**Logit Attribution** We use logit attribution to examine specific information (e.g., operator or answer tokens) at each layer (See Figure 4a for illustration). Specifically, we compute *direct logit attribution* (nostalgebraist, 2020; Belrose et al., 2023) of a target token  $t$  by projecting hidden states at various points in each layer onto the vocabulary space:  $\text{logit}(t) = \langle W_U[t], \text{LN}(h) \rangle$ , where  $h$  is the hidden state,  $\text{LN}$  is LayerNorm, and  $W_U[t]$  is the unembedding vector. We probe four points within each layer at the last token position: the attention output, MLP output, and the residual stream immediately after merging the attention output (*resid mid*) and after merging the MLP output (*resid final*). As summarized in Figure 4a, we track abstraction via the logits of operator tokens (e.g., “+”, “add”, “addition”) and computation via the logits of operand and answer tokens across layers.

---

#### Algorithm 1 Activation Patching

---

- 1: **Input:** Set  $\Omega$  of clean and corrupted sample pairs  $(X_{cl}, X_{cor})$ , model  $\mathcal{M}$  with hidden states  $\mathcal{S}$ .
  - 2: **Output:** Patching effects for  $\mathcal{S}$ :  $E_{\mathcal{S}}$ .
  - 3: **for**  $(X_{cl}^{(i)}, X_{cor}^{(i)}) \in \Omega$  **do**
  - 4:    $\text{logit}_o, A_{cl} \leftarrow \mathcal{M}(X_{cl}^{(i)}, A_{cl})$  # Clean run: get clean logits and store all layer activations  $A_{cl}$
  - 5:    $\text{logit}_c, A_{cor} \leftarrow \mathcal{M}(X_{cor}^{(i)}, A_{cor})$  # Corrupted run: get corrupted logits and store all layer activations  $A_{cor}$
  - 6:   **for**  $s \in \mathcal{S}$  **do**
  - 7:      $A'_{cor}(s) \leftarrow A_{cl}(s)$  # Patched run: replace hidden state  $s$  in  $A_{cor}$  by  $A_{cl}$
  - 8:      $\text{logit}_p \leftarrow \mathcal{M}(X_{cor}^{(i)}, A'_{cor})$  # get patched logits
  - 9:      $e_s^{(i)} \leftarrow \frac{\text{logit}_p - \text{logit}_c}{\text{logit}_o - \text{logit}_c}$  # patching effect
  - 10:   **end for**
  - 11: **end for**
  - 12: **Return:**  $E_{\mathcal{S}} \leftarrow \frac{1}{|\Omega|} \sum_{i=1}^{|\Omega|} e_s^{(i)}$
- 

**Activation Patching** To identify components causally responsible for abstraction and computation, we apply activation patching (Algorithm 1, see Figure 4b for visualization) (Ghandeharioun et al., 2024; Zhang and Nanda, 2023; Meng et al., 2022). To quantify the contribution of each component across layers, this method replaces a single

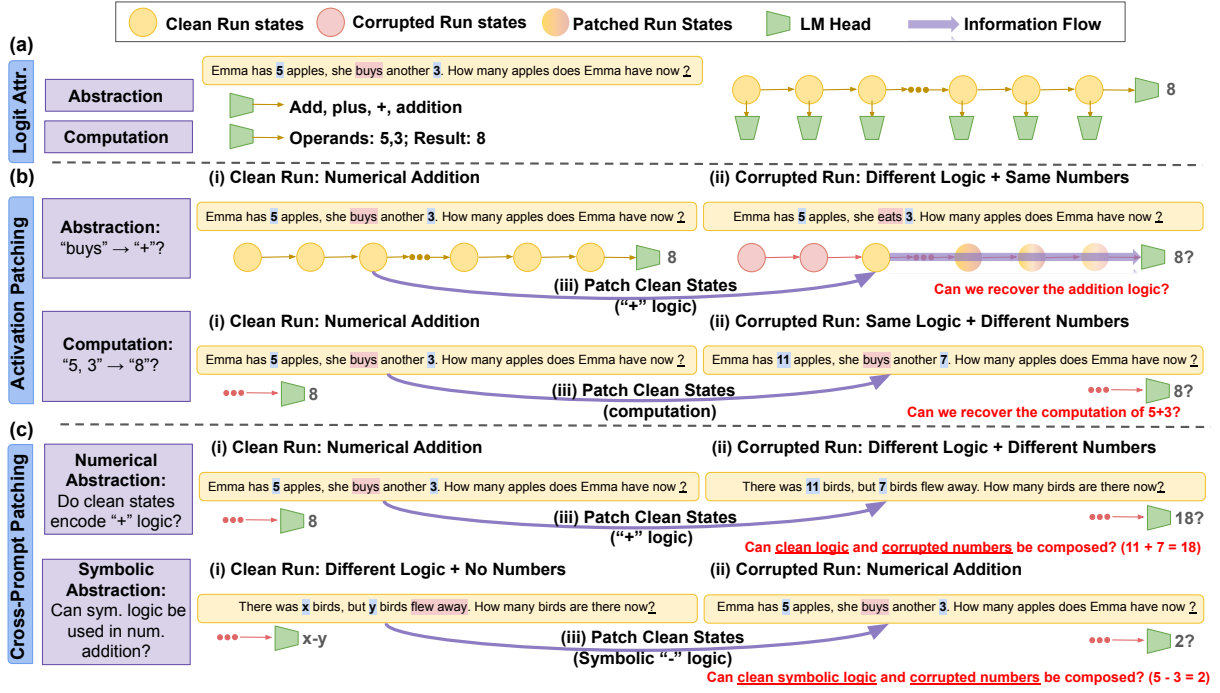


Figure 4: Overview of interpretability methods probing the abstract-then-compute mechanism in simple math problems, focusing on hidden states at the last token position across layers.

intermediate hidden state in the corrupted forward pass with the corresponding hidden state from the clean run and measures how much this single hidden state injected in corrupted forward pass can restore the prediction of the clean answer. This patching effect per state per layer is a normalized score from 0 (no recovery) to 1 (full recovery to clean performance), with higher indicating more contribution. We patch attention, MLP and final layer outputs across layers at the last position. Formally, we quantify causal impact using the logit difference between clean  $a_{cl}^{(i)}$  and corrupted answers  $a_{cor}^{(i)}$  in Eq. 2.

$$LD_*(i) = \text{logit}_*(a_{cl}^{(i)}) - \text{logit}_*(a_{cor}^{(i)}) \quad (1)$$

$$e_s^{(i)} = \frac{LD_p(i) - LD_c(i)}{LD_o(i) - LD_c(i)} \quad (2)$$

To probe **abstraction** (Figure 4b), we construct minimally different clean/corrupted pairs that vary in their underlying logic (e.g., “buys” for addition vs. “eats” for subtraction) but have the same numbers (e.g., “5,3”). In Figure 4b, the clean input implies  $5 + 3 = 8$ , while the corrupted input implies  $5 - 3 = 2$ . We patch individual clean states to the corrupted run to identify critical layers for restoring the addition logic and recovering the clean answer ‘8’. For **computation** (Figure 4b), we use pairs with the same logic (e.g., addition), but different

numbers (e.g., “5,3” vs. “11,7”). Here, we seek to identify layers whose states when patched individually from the clean run to the corrupted run, are critical to perform the clean-run-specific computation with clean operands 5, 3 and output “8”.

### 5.1.2 Abstract-then-Compute Hypothesis

As shown in Figure 5, we observe distinct stages for abstraction and computation, supporting the *abstract-then-compute* hypothesis. Logit attribution reveals that around L13–14, attention begins moving the inferred operator (e.g., ‘+’, plus’, add’) to the last position (Figure 5i, iv). This coincides with a divergence in logit differences between target operators (+’ vs. ‘-’) in addition and subtraction problems (Figure 5v), suggesting that while earlier layers encode generic operator features, problem-specific abstraction emerges here. Subsequently, around L15–16, Figure 5i,iv shows operands transfer to the last position; Following abstraction, the computation phase appears to begin at L18, primarily through MLPs layers (Figure 5ii, iv). Activation patching confirms the distinct stages: abstraction starts at around L13, with rising attention and layer patching effects (Figure 5iii); The rise of attention and layer patching effects in L15,16 in Figure 5vi aligns with our previous observation that operands are being moved to the last position. Finally, the peak patching effect of MLP at L18 high-

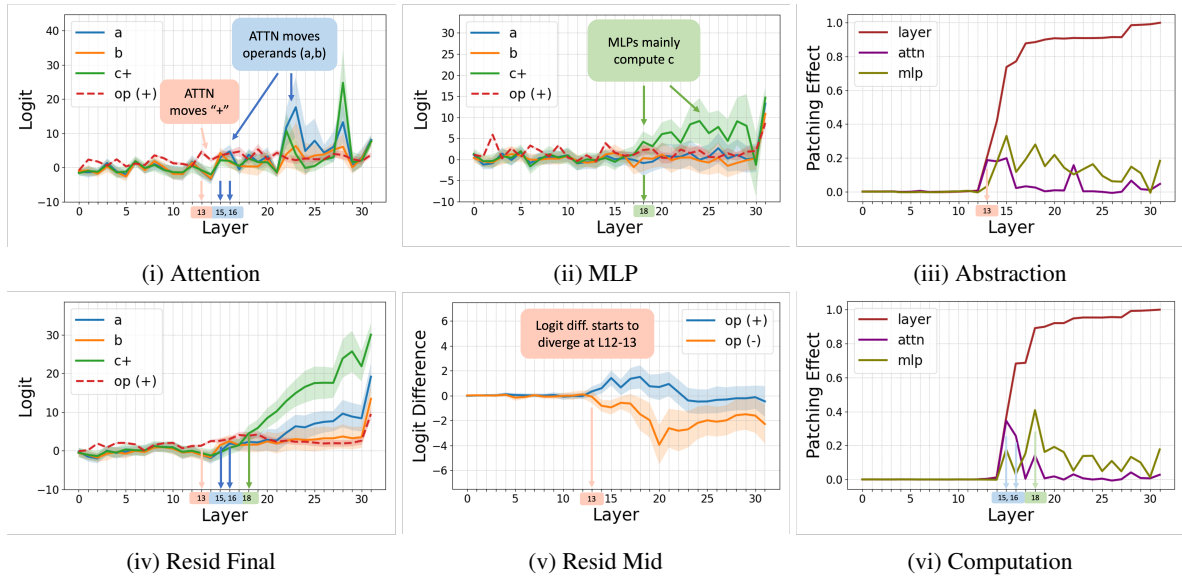


Figure 5: Visualizations of internal computations at last token position in Llama-3 8B for addition math word problems: (i,ii, iv, v) for logit attribution results where  $a, b$  are operands and  $c$  is the result; (iii, vi) for activation patching results. We label the starting layer of abstraction, operand moving and computation in pink, blue and green, respectively.

light their crucial role in calculating the answer. These combined results support our hypothesis that the model follows an *abstract-then-compute* mechanism within a single forward pass. Additional logit attribution and activation patching results for other models and two-operator problems are in Appendix B.2.

## 5.2 Validation and Abstraction Transfer with Cross-Prompt Patching

We now validate the causal role of the critical layers for abstraction (L13,14) and computation (L15,16 for operands, L18 for execution). We also investigate if the abstraction representations formed at around L13,14 are transferable across problem forms (symbolic/concrete) and templates, and can be composed with subsequent computation stage.

**Method** Cross-prompt patching also uses Algorithm 1, but instead of computing patching effects, we track the log-probability of specific tokens across layers in each patched run. This acts as a form of “knock-out” intervention: we overwrite a single layer’s activations in the corrupted run with clean activations that are hypothesized to contain specific information (e.g., abstraction, operands, computation), and observe whether this information is reflected in the output level.

To validate the critical layers for each stage, we cross-patch for **numerical abstraction** (Figure 4c), where both the clean and corrupted inputs are nu-

merical problems, but differ in both underlying logic and operands. As shown in Figure 4c, the clean run corresponds to  $5 + 3 = 8$  and the corrupted to  $11 - 7 = 4$ . We patch hidden states from the clean run into the corrupted run and validate our hypothesis: (i) *Abstraction (L13-14)*: At these layers, operands have not been transferred yet, so patching should only transfer the clean logic. If these layers encode addition logic, the model should apply the clean addition operator to the corrupted operands, computing  $11 + 7 = 18$  in the remaining forward pass. We expect the log-probability of this *target answer* (‘18’) to rise. (ii) *Operand Transfer (L15-16)*: L15 begins operand transfer and already contains both clean run logic and operand information. Patching them should increase the log-prob of the clean answer (e.g., ‘8’), while reducing probability of the corrupted (‘4’) and target answers (‘18’). (iii) *Computation (Layer 18)*: By this point, the full ingredients (abstraction and computation) are available. Patching here should fully recover the clean answer (‘8’). We expect the log-prob close to 0. To evaluate these effects, we track log-probabilities across layers for the **target answer** ( $18 = 11 + 7$ , clean logic + corrupted operands) – testing numerical abstraction transfer, the **clean answer** ( $8 = 5 + 3$ ) – testing operand alignment and execution, and the **corrupted answer** ( $4 = 11 - 7$ ).

To investigate if the abstraction representations

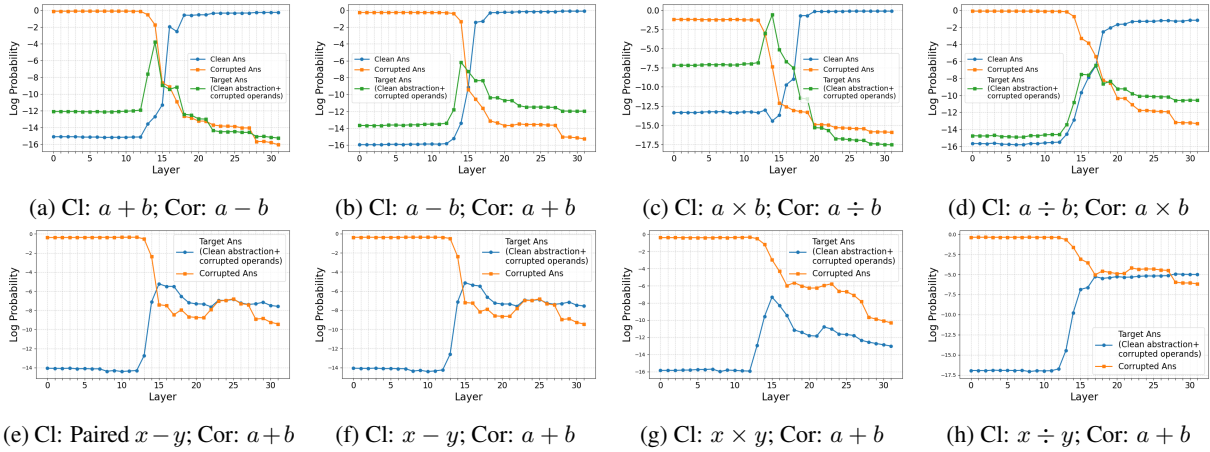


Figure 6: Cross-patching results for Llama-3 8B with corresponding clean and corrupted run.  $a, b$  indicate concrete numerical problems, while  $x, y$  indicate symbolic problems. **Top (Numerical Abstraction):** Patching concrete problems with different abstractions shows target log-prob rising at 13 (abstraction onset), peaking at 14 (abstraction formed), then falling as clean operands are introduced. Meanwhile, the clean answer’s log-prob rises from 13 (abstraction) and 15 (operand integration), stabilizing at layer 18 (computation). **Bottom (Symbolic Abstraction):** Patching symbolic problems into concrete addition shows target log-probability rising at layer 13, peaking at 15 (where predictions flip), then declining.

can be transferred across problem forms (symbolic/numerical) and templates, and if they can be composed with subsequent computation stage, we cross-patch for **symbolic abstraction** (See Figure 4c) – patching symbolic clean states to numerical corrupted run. Here, clean inputs are symbolic (no concrete numbers), and corrupted inputs are numerical problems with a different underlying logic. This ensures that only abstraction (no operands or computation) is transferred from the clean run, unlike numerical abstraction cross-patching. In the example in Figure 4c, clean run predicts  $x - y$ , while the corrupted run corresponds to  $5 + 3 = 8$ . By patching clean states from the symbolic problem to the numerical corrupted run, we examine (i) if *symbolic abstractions are also formed* at around L13-14, despite predicting ‘x’ as the first token, and (ii) if this abstraction ( $x - y$ ), when transferred into numerical corrupted run, can be *used and composed* with corrupted operands (5, 3) to compute  $5 - 3 = 2$ . To assess this, we track the per-layer log-probabilities of the **target answer** (clean logic + corrupted operands,  $2 = 5 - 3$ ) and **corrupted answer** ( $8 = 5 + 3$ ), and omit the clean answer ‘x’. If symbolic abstraction transfer occurs, we expect an increase in the target answer log-prob, and a corresponding decrease in the corrupted answer starting around L13. Note that since the symbolic clean states across layers are predicting ‘x’, we expect both answer log-probs to drop.

**Results** Figure 6a shows results for **numerical abstraction** cross-patching results corresponding to the example illustrated in Figure 4. As expected, the target answer log-probability (‘18’) begins rising at L13 (abstraction onset), peaks at L14 (abstraction formed), and drops when clean operands are introduced (L15). The clean answer (‘8’) log-probability keeps rising from L13 (abstraction) and continue at 15 (operand integration), stabilizing by L18 (computation). The corrupted answer (‘4’) log-probability drops after L13. These trends hold across underlying logic (Figure 6b-d), confirming the roles of these critical layers as identified earlier. In **symbolic abstraction** cross-patching (Figure 6e-h), we observe consistent behaviour: from L13 onward, the target answer probability increases while the clean answer decreases, eventually flipping. This indicates that (i) abstractions injected via patching are composed with corrupted operands to produce valid outputs, and (ii) abstraction representations at L13–14 are invariant to surface form and problem template. Concretely, comparing Figure 6e and Figure 6f, where minimally different templates are used in (e) and random templates in (f), we observe near-identical effects in both cases —suggesting abstraction transfer is template-invariant. Furthermore, (g) and (h) show that injecting symbolic *multiplication* and *division* abstractions into concrete *addition* problems still flips the model’s prediction—demonstrating the generality of abstraction transfer. Cross-patching results for



other models, and two-operator problems are in Appendix B.3.

Together, these results provide strong support for the abstract-then-compute hypothesis with critical layers for abstraction (L13,14) and computation (L15,16 for operands and L18 for computation), and further demonstrate that: (i) abstraction can be transferred and composed with subsequent computation across surface forms (symbolic/concrete) and templates, and (ii) even at the last position in symbolic problems, when predicting the first output token ‘x’, middle layers already encode abstraction (e.g., the correct operator), indicating that next-token prediction reflects not just immediate token prediction, but also anticipates future outputs.

## 6 Discussion

Math reasoning tasks, particularly structured and easily verifiable math word problems such as GSM8K, have attracted significant attention in the community. Our findings provide a complementary perspective to recent work on efficient and latent reasoning.

**Toward More Token-Efficient CoT Reasoning:** Recent work has explored reducing the token cost of lengthy CoT chains (Ma et al., 2025; Fan et al., 2025; Fatemi et al., 2025; Munkhbat et al., 2025; Kang et al., 2025; Xia et al., 2025). Our findings offer a complementary perspective: if models internally represent abstract problem structures and the computation step is the primary bottleneck, then the allocation of intermediate tokens could be reconsidered. While CoT facilitates arithmetic computation, it raises a key question: what is the minimal number of CoT steps required for accurate arithmetic computation? This suggests new directions for designing concise yet effective CoT scaffolds that exploit internal abstraction, reducing verbosity without sacrificing performance.

**Latent CoT and Internal Computation:** Our results suggest that the abstraction skills exist internally, which supports the emerging view that models can internalize reasoning, even when such reasoning steps are not explicitly verbalized in the generated output. This is consistent with recent work on latent CoT (Hao et al., 2024; Pfau et al., 2024; Goyal et al., 2024; Cheng and Van Durme, 2024; Su et al., 2025; Zhang et al., 2025; Wu et al., 2025), where the reasoning is embedded in the hidden space rather than in discrete vocabulary space.

The abstract-then-compute behavior and abstraction transferability across problems suggest that reasoning may be increasingly internalized with scale and training. These insights can inform future work on dynamic strategy selection during decoding—where models might conditionally decide whether to externalize reasoning or rely on internal computation—depending on task difficulty, context, or resource constraints.

## 7 Conclusion

Disentangled evaluation reveals that, without CoT, models perform better at abstraction than computation, with the latter bottlenecking final-answer accuracy — challenging the view that poor performance always implies reasoning failure. Mechanistic interpretability uncovers an abstract-then-compute mechanism with transferable abstractions. We argue for disentangled evaluation to more precisely assess model abilities and inform architectural design.

## 8 Limitations

Our study has several limitations. First, we focus solely on English-language datasets; whether the abstract-then-compute mechanism generalizes to other languages remains an open question. Second, our evaluation decomposes mathematical problem-solving into only two stages: abstract formulation and arithmetic computation. Finer-grained breakdowns (e.g., Opedal et al. (2024)) may offer deeper insight. Third, our interpretability analysis is limited to single-step generation, as common techniques (e.g., activation patching, logit attribution) target single-token behavior. Extending these to multi-step reasoning is an ongoing challenge, with recent work like SelfIE (Chen et al., 2024) provides initial steps. Fourth, while we focus on critical layers involved in abstraction and computation, we leave detailed analysis of components to future work. Finally, due to compute constraints, we analyze models up to 12B parameters. Extending to larger models is left for future studies.

## Acknowledgments

This work is supported by the Mitacs Accelerate Program (Project ID: IT27067) and partially enabled by Mila’s computing resources (mila.quebec). We thank Zichao Li, Cesare Spinoso-Di Piano, and Xiyuan Zou for helpful discussions. Jackie Chi Kit Cheung is supported by the Canada CIFAR

AI Chair program. We acknowledge the material support of NVIDIA for providing computational resources.

## References

- Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. [MathQA: Towards interpretable math word problem solving with operation-based formalisms](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics.
- Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. 2023. Eliciting latent predictions from transformers with the tuned lens. *arXiv preprint arXiv:2303.08112*.
- Haozhe Chen, Carl Vondrick, and Chengzhi Mao. 2024. [SelfIE: Self-interpretation of large language model embeddings](#). In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 7373–7388. PMLR.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.
- Jeffrey Cheng and Benjamin Van Durme. 2024. Compressed chain of thought: Efficient reasoning through dense representations. *arXiv preprint arXiv:2412.13171*.
- Ziling Cheng, Meng Cao, Marc-Antoine Rondeau, and Jackie Chi Kit Cheung. 2025. [Stochastic chameleons: Irrelevant context hallucinations reveal class-based \(mis\)generalization in llms](#). *Preprint*, arXiv:2505.22630.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Chenrui Fan, Ming Li, Lichao Sun, and Tianyi Zhou. 2025. Missing premise exacerbates overthinking: Are reasoning models losing critical thinking skill? *arXiv preprint arXiv:2504.06514*.
- Mehdi Fatemi, Banafsheh Rafiee, Mingjie Tang, and Kartik Talamadupula. 2025. Concise reasoning via reinforcement learning. *arXiv preprint arXiv:2504.05185*.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Pal: Program-aided language models. In *International Conference on Machine Learning*, pages 10764–10799. PMLR.
- Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lucas Dixon, and Mor Geva. 2024. Patchscopes: A unifying framework for inspecting hidden representations of language models. *arXiv preprint arXiv:2401.06102*.
- Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. 2023. Localizing model behavior with path patching. *arXiv preprint arXiv:2304.05969*.
- Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. 2024. [Think before you speak: Training language models with pause tokens](#). In *The Twelfth International Conference on Learning Representations*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. 2024. Training large language models to reason in a continuous latent space. *arXiv preprint arXiv:2412.06769*.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. 2024. [OlympiadBench: A challenging benchmark for promoting AGI with olympiad-level bilingual multimodal scientific problems](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3828–3850, Bangkok, Thailand. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021a. [Measuring massive multitask language understanding](#). In *International Conference on Learning Representations*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. Measuring mathematical problem solving with the math dataset. *NeurIPS*.
- Yu Kang, Xianghui Sun, Liangyu Chen, and Wei Zou. 2025. C3ot: Generating shorter chain-of-thought without compromising effectiveness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 24312–24320.
- Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. 2016. [MAWPS: A math word problem repository](#). In *Proceedings of*

- the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: *Human Language Technologies*, pages 1152–1157, San Diego, California. Association for Computational Linguistics.
- Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. 2024. Gsm-plus: A comprehensive benchmark for evaluating the robustness of llms as mathematical problem solvers. *arXiv preprint arXiv:2402.19255*.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 158–167, Vancouver, Canada. Association for Computational Linguistics.
- Yu Lu Liu, Su Lin Blodgett, Jackie Cheung, Q. Vera Liao, Alexandra Olteanu, and Ziang Xiao. 2024. ECBD: Evidence-centered benchmark design for NLP. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16349–16365, Bangkok, Thailand. Association for Computational Linguistics.
- Wenjie Ma, Jingxuan He, Charlie Snell, Tyler Griggs, Sewon Min, and Matei Zaharia. 2025. Reasoning models can be effective without thinking. *arXiv preprint arXiv:2504.09858*.
- Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. In *Advances in Neural Information Processing Systems*.
- Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. 2024. Language models implement simple Word2Vec-style vector arithmetic. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5030–5047, Mexico City, Mexico. Association for Computational Linguistics.
- 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.
- Seyed Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2025. GSM-symbolic: Understanding the limitations of mathematical reasoning in large language models. In *The Thirteenth International Conference on Learning Representations*.
- Swaroop Mishra, Matthew Finlayson, Pan Lu, Leonard Tang, Sean Welleck, Chitta Baral, Tanmay Rajpurohit, Oyvind Taffjord, Ashish Sabharwal, Peter Clark, and Ashwin Kalyan. 2022. LILA: A unified benchmark for mathematical reasoning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5807–5832, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Tergel Munkhbat, Namgyu Ho, Seo Hyun Kim, Yongjin Yang, Yujin Kim, and Se-Young Yun. 2025. Self-training elicits concise reasoning in large language models. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 25127–25152, Vienna, Austria. Association for Computational Linguistics.
- Yaniv Nikankin, Anja Reusch, Aaron Mueller, and Yonatan Belinkov. 2025. Arithmetic without algorithms: Language models solve math with a bag of heuristics. In *The Thirteenth International Conference on Learning Representations*.
- nostalgebraist. 2020. interpreting gpt: the logit lens.
- Andreas Opedal, Alessandro Stolfo, Haruki Shirakami, Ying Jiao, Ryan Cotterell, Bernhard Schölkopf, Abulhair Saparov, and Mrinmaya Sachan. 2024. Do language models exhibit the same cognitive biases in problem solving as human learners? *arXiv preprint arXiv:2401.18070*.
- OpenAI. 2024. Gpt-4o mini: Advancing cost-efficient intelligence. <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>.
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP models really able to solve simple math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, Online. Association for Computational Linguistics.
- Jacob Pfau, William Merrill, and Samuel R. Bowman. 2024. Let’s think dot by dot: Hidden computation in transformer language models. In *First Conference on Language Modeling*.
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In *International Conference on Machine Learning*, pages 31210–31227. PMLR.
- Shuming Shi, Yuehui Wang, Chin-Yew Lin, Xiaojiang Liu, and Yong Rui. 2015. Automatically solving number word problems by semantic parsing and reasoning. In *Proceedings of the 2015 conference on empirical methods in natural language processing*, pages 1132–1142.
- Zayne Rea Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. 2025. To cot or not to cot? chain-of-thought helps mainly on math and symbolic reasoning. In *The Thirteenth International Conference on Learning Representations*.

DiJia Su, Hanlin Zhu, Yingchen Xu, Jiantao Jiao, Yuandong Tian, and Qinqing Zheng. 2025. Token assorted: Mixing latent and text tokens for improved language model reasoning. *arXiv preprint arXiv:2502.03275*.

Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2023. [Interpretability in the wild: a circuit for indirect object identification in GPT-2 small](#). In *The Eleventh International Conference on Learning Representations*.

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*.

Junhong Wu, Jinliang Lu, Zixuan Ren, Ganqiang Hu, Zhi Wu, Dai Dai, and Hua Wu. 2025. Llms have a heart of stone: Demystifying the soft thinking ability of large reasoning models. *arXiv preprint arXiv:2508.03440*.

Heming Xia, Chak Tou Leong, Wenjie Wang, Yongqi Li, and Wenjie Li. 2025. Tokenskip: Controllable chain-of-thought compression in llms. *arXiv preprint arXiv:2502.12067*.

Qwen An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Junyang Lin, and 25 others. 2024. [Qwen2.5 technical report](#). *ArXiv*, abs/2412.15115.

Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-Zhu. 2025. [Physics of language models: Part 2.1, grade-school math and the hidden reasoning process](#). In *The Thirteenth International Conference on Learning Representations*.

Zeping Yu and Sophia Ananiadou. 2024. [Interpreting arithmetic mechanism in large language models through comparative neuron analysis](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 3293–3306, Miami, Florida, USA. Association for Computational Linguistics.

Fred Zhang and Neel Nanda. 2023. Towards best practices of activation patching in language models: Metrics and methods. *arXiv preprint arXiv:2309.16042*.

Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, William Song, Tiffany Zhao, Pranav Vishnu Raja, Charlotte Zhuang, Dylan Z Slack, Qin Lyu, Sean M. Hendryx, Russell Kaplan, Michele Lunati, and Summer Yue. 2024a. [A careful examination of large language model performance on grade school arithmetic](#). In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Wei Zhang, Chaoqun Wan, Yonggang Zhang, Yiu-ming Cheung, Xinmei Tian, Xu Shen, and Jieping Ye. 2024b. Interpreting and improving large language models in arithmetic calculation. *arXiv preprint arXiv:2409.01659*.

Xiaotian Zhang, Chunyang Li, Yi Zong, Zhengyu Ying, Liang He, and Xipeng Qiu. 2023. Evaluating the performance of large language models on gaokao benchmark. *arXiv preprint arXiv:2305.12474*.

Zhen Zhang, Xuehai He, Weixiang Yan, Ao Shen, Chenyang Zhao, Shuohang Wang, Yelong Shen, and Xin Eric Wang. 2025. Soft thinking: Unlocking the reasoning potential of llms in continuous concept space. *arXiv preprint arXiv:2505.15778*.

Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2024. [AGIEval: A human-centric benchmark for evaluating foundation models](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2299–2314, Mexico City, Mexico. Association for Computational Linguistics.

Tianyi Zhou, Deqing Fu, Vatsal Sharan, and Robin Jia. 2024. Pre-trained large language models use fourier features to compute addition. *Advances in Neural Information Processing Systems*, 37:25120–25151.

## A Disentangled Evaluation Details and Additional Results

### A.1 Symbolic Variant Creation For GSM8K and SVAMP

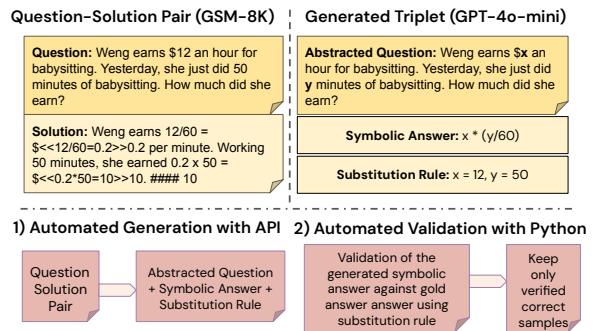


Figure 7: **Generate-then-validate pipeline:** We use API calls to obtain abstract question-answer-substitution triplets from the concrete question-solution pair from GSM-8K, then validate them against gold answer using sympy. Triplets that fail this check are manually reviewed and corrected.

All our evaluations are conducted on the GSM8K test set and the full SVAMP dataset. To support our evaluation of abstract formulation and arithmetic computation in Section 4, we construct symbolic question and expression answer variants for both

SVAMP and GSM-8K. Examples are shown in Table 4.

For SVAMP, which already includes both the expression (e.g.,  $20 \times 10$ ) and the final numerical answer (e.g., 200), we create symbolic abstraction variants by replacing all numeric values with symbolic variables (e.g.,  $x, y$ ) in both the question and the corresponding expression. This preserves the structure and semantics of the original problem while abstracting away from the concrete numbers. For arithmetic computation variant, we use the paired expression.

For GSM-8K, which lacks such annotations, we generate both the symbolic abstraction variant and the numerical expressions using a two-stage generate-then-validate pipeline (Figure 7). In the generation stage, we use GPT-4o-mini (OpenAI, 2024) to produce triplets from original question–solution pairs. Each triplet consists of: (1) a symbolic version of the question, where relevant numbers are replaced with variables while maintaining the semantic content; (2) a symbolic expression that represents the solution in closed-form using those variables; and (3) a substitution rule that maps each variable to its original numeric value. In the validation stage, we verify the correctness of each generated sample. We apply the substitution rule to the symbolic expression, obtaining a numerical expression, then using `sympy` to evaluate the expression, and compare the resulting numeric answer to the gold answer from GSM-8K. Triplets that fail this check are manually reviewed and corrected.

## A.2 Evaluation Details

In this section, we detail the evaluation of the four settings. First, we show the instructions used in each settings in Table 5 with and without CoT. The prompt used in each setting is then a concatenation of the instruction and the question.

For the *Original* and *Arithmetic Computation* settings, where the expected output is a final integer answer, we extract the answer following the token #####, remove any accompanying units, and normalize formatting (e.g., removing commas, dollar signs, percentage symbols, and units like ‘g’) before comparing it with the gold answer.

For the *Numerical Abstraction* setting, where answers are expected to be *numerical expressions*, we first convert LaTeX-style expressions to Python syntax (when written in Markdown form), then evaluate them using `sympy` to check equivalence

### Box 1: Symbolic Evaluation Prompt

Determine whether the following two mathematical expressions are equivalent. The expressions may be written in simplified or unsimplified symbolic form (e.g.,  $1/2x + 3$ ), natural language (e.g., “Susan made  $1/2x + 3$  buttons”) or in LaTeX notation. Consider expressions equivalent if they represent the same mathematical value, even if written differently (e.g., different notation, simplification, or variable order when valid). Respond only with: True or False.

#### Example:

1.  $z - (y - x)$

2.  $z - y + x$

**Answer:** True

1. Susan made  $1/2 \times x$  buttons

2.  $0.5x$

**Answer:** True

1.  $2(y + x)$

2.  $M = 2(y + x)$

**Answer:** True

1.  $xz * ((1 - y)/100)$

2.  $x * z - (y/100) * (x * z)$

**Answer:** True

#### Now evaluate:

1. `{symbolic_gold_answer}`

2. `{abstract_generated_answer}`

**Answer:**

with the gold expression.

In the *Symbolic Abstraction* setting, where outputs are *symbolic expressions*, we use gpt-4o-mini as an automated evaluator. The prompting to gpt-4o-mini is shown in Box 1, and responses are generated with temperature set to 0. To validate this method, we annotated a held-out set of 120 samples manually for correctness, and compared our annotations with the gpt-4o-mini evaluator’s decisions. We find that gpt-4o-mini achieves **94% agreement** with our judgment in identifying symbolic expression equivalence. Example comparisons are shown in Table 3.

## A.3 Additional Result of Disentangled Evaluation Without CoT

We report zero-shot, no-CoT performance on SVAMP in Figure 9. Compared to GSM8K, SVAMP is a significantly simpler benchmark consisting of math word problems that require only a single reasoning step — namely, a single arithmetic operation. As with GSM8K, models perform better on the abstraction variants than in the original setting, though the performance gap is smaller due

Gold Answer	Model Generation	Our Eval	GPT-4o-mini Eval
$u * (x + y + z)$	$xu + yu + zu$	True	True
$x + x * (1/y)$	$x + (x/y)$	True	True
$0.5(x + yz)$	$z * (y + 1) * x/2$	False	False
$(y + z)/x$	$xz - y = xy$	False	False
$xz * ((1 - y)/100)$	$(x * (1 - y/100) * z)$	False	True
$(12/x) * y$	$y * 12$	False	True

Table 3: Comparison of gold answers, model generations, our annotated correctness, and GPT-4o-mini evaluation on a held-out set of 120 samples.

Dataset	Symbolic Question	Answer	Substitution
GSM8K	I have $x$ liters of orange drink that are $y\%$ water and I wish to add it to $z$ liters of pineapple drink that is $u\%$ water. But as I pour it, I spill $v$ liters of the orange drink. How much water is in the remaining $w$ liters?	$(y \cdot (x - v) + u \cdot z) / 100$	$x = 10, y = \frac{2}{3}, z = 15, u = \frac{3}{5}, v = 1, w = 24$
GSM8K	Jerry has a flock of chickens. The red chickens produce $x$ eggs a day, and the white chickens produce $y$ eggs a day. Every day Jerry collects $z$ eggs. If he has $u$ more white chickens than red chickens, how many red chickens does he have?	$(z - u \cdot y) / (x + y)$	$x = 3, y = 5, z = 42, u = 2$
GSM8K	Adrian’s age is $x$ times the age of Harriet, and Harriet is $y$ the age of Zack. Calculate the average age of the three in three years if Harriet is $z$ years old now.	$(x * z + z + (z/y) + 9) / 3$	$x = 3, y = \frac{1}{2}, z = 21$
SVAMP	Each pack of DVDs costs $x$ dollars. If there is a discount of $y$ dollars on each pack	$x - y$	$x = 76, y = 25$
SVAMP	An industrial machine worked for $x$ minutes. It can make $y$ shirts a minute.	$x \cdot y$	$x = 4, y = 5$
SVAMP	Paco had $x$ salty cookies and $y$ sweet cookies. He ate $z$ sweet cookies and $u$ salty cookies. How many salty cookies did Paco have left?	$x - u$	$x = 26, y = 17, z = 14, u = 9$

Table 4: Constructed symbolic examples from GSM8K and SVAMP datasets.

to the task’s simplicity.

Interestingly, we observe a notable difference from GSM8K: across all model sizes, even small models such as LLAMA 1B and 3B perform well on the *Arithmetic Computation* variant, often outperforming both the abstraction variants and the original setting. This suggests that computing one-step expressions (e.g.,  $5 - 3$ ) is less challenging than deriving an abstract formulation with only one step. However, in tasks involving multiple steps, abstraction becomes comparatively easier than executing the full computation correctly, as shown in the case of GSM8K. This highlights how model capabilities depend not just on the skill type but also on the complexity of the required operation.

#### A.4 Additional Results of Disentangled Evaluation With CoT

We present the full results on GSM8K for Llama family and Qwen family in Figure 10, and full results on SVAMP for Llama family and Qwen

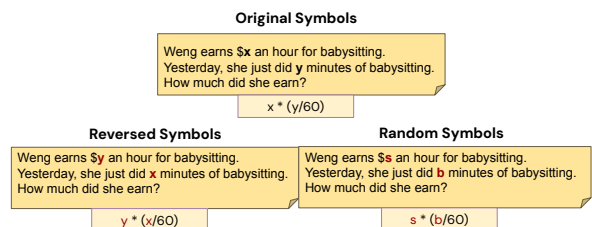


Figure 8: Experiment configurations for the ablation study on symbol choices and symbol order.

family in Figure 11.

#### A.5 Ablation on Symbolic Abstraction Variant

To assess the reliability and external validity of the symbolic abstraction evaluation, we perform ablations over symbol order and symbol choice. As illustrated in Figure 8, we compare three settings:

- **Original Symbols:** Variables are consistently represented using a fixed set of letters in order— $x, y, z, u, v, w, p, q, r, s, t$ —e.g.,

Table 5: Prompting Strategies, Problem Variants and Instructions

Setting	Strategy	Instruction	Question	Answer
Original	No CoT	Please answer the question directly WITHOUT showing the reasoning process, you MUST write the answer as <b>an integer</b> after '####', without including the equation or units.	Weng earns \$12 an hour for babysitting. Yesterday, she just did 50 minutes of babysitting. How much did she earn?	10
Original	CoT	<b>Let's think step by step</b> , you MUST write the answer as <b>an integer</b> after '####' without including the units. Write the answer at the end.	Weng earns \$12 an hour for babysitting. Yesterday, she just did 50 minutes of babysitting. How much did she earn?	10
Arithmetic Computation	No CoT	Please answer the question directly WITHOUT showing the reasoning process, you MUST write the answer as <b>an integer</b> after '####'	What is the value of $12 * (50/60)$ ?	10
Arithmetic Computation	CoT	<b>Let's think step by step</b> , you MUST write the answer as <b>an integer</b> after '####' . Write the answer at the end.	What is the value of $12 * (50/60)$ ?	10
Numerical Abstraction	No CoT	Please answer the question directly without showing the reasoning process, you MUST write the <b>expression</b> with appropriate round brackets after '####', without including the units, and you DO NOT need to simplify the expression.	Weng earns \$12 an hour for babysitting. Yesterday, she just did 50 minutes of babysitting. How much did she earn?	$12 * (50/60)$
Numerical Abstraction	CoT	<b>Let's think step by step</b> , at the end, you MUST write the <b>expression</b> with appropriate parenthesis after '####', without including the units, but you DO NOT need to simplify the expression.	Weng earns \$12 an hour for babysitting. Yesterday, she just did 50 minutes of babysitting. How much did she earn?	$12 * (50/60)$
Symbolic Abstraction	No CoT	Please answer the question directly WITHOUT showing the reasoning process, you MUST write the <b>expression</b> with appropriate round brackets after '####' without including the units, and you DO NOT need to simplify the expression.	Weng earns \$x an hour for babysitting. Yesterday, she just did y minutes of babysitting. How much did she earn?	$x * (y/60)$
Symbolic Abstraction	CoT	<b>Let's think step by step</b> , at the end, you MUST write the <b>expression</b> with appropriate round brackets after '####' without including the units, but you DO NOT need to simplify the expression.	Weng earns \$x an hour for babysitting. Yesterday, she just did y minutes of babysitting. How much did she earn?	$x * (y/60)$

$$x \times (y/60).$$

- **Reversed Symbols:** The same set of symbols is used, but the order is reversed (e.g.,  $y \times (x/60)$ ), preserving the semantic and structural content of the problem while changing the superficial presentation.
- **Random Symbols:** Each original symbol is replaced with a randomly sampled letter from the alphabet, unique to each dataset. This preserves the structure of the expression while removing any consistent identity cues. The mappings are as follows: {'a': 'h', 'd': 'i', 'm': 's', 'n': 'r', 'p': 'e', 'q': 'l', 'r': 'c', 's': 'v', 't': 'j', 'u': 'm', 'v': 't', 'w': 'o', 'x': 'u', 'y': 'p', 'z': 'b', 'Z': 'f'}

In Table 6, we observe mild performance degradation with symbol perturbations on both models, (e.g., three-point drop with Reversed and another two points with Random), but models retain strong accuracy compared to the Original setting. This

suggests that Symbolic Abstraction is relatively robust to surface-level symbol changes.

Setting	Llama 8B		Qwen 7B	
	No CoT	CoT	No CoT	CoT
original	45.7	56.7	61.5	74.7
reverse	42.8	51.8	61.9	74.8
random	41.0	53.1	58.0	71.9

Table 6: Results of ablation study on symbol choices and symbol order, with and without CoT under zero-shot setting on GSM8K.

## B Interpretability Results

### B.1 Interpretability Data Construction

To construct a dataset suitable for mechanistic interpretability, we focus on simpler math word problems that require only one or two reasoning steps with one or two basic arithmetic operations (addition, subtraction, multiplication, or division). We deliberately avoid more complex multi-step problems, as model performance on such tasks tends to be poor, potentially confounding interpretability

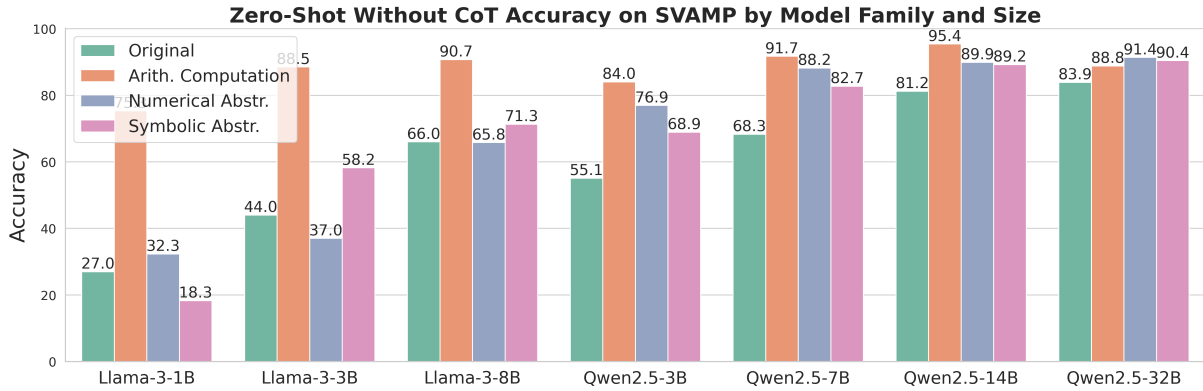


Figure 9: Model zero-shot **without** CoT performance on SVAMP.



Figure 10: Model zero-shot **with and without** CoT performance on GSM8K. A.C.: Arithmetic Computation; N.A.: Numerical Abstraction; O.: Original; S.A.: Symbolic Abstraction.

analyses.

For each pair of arithmetic operations— $(x + y, x - y)$  and  $(x \times y, x \div y)$  and  $(x + y + z, x + y - z, x - y + z, x - y - z)$ —we use a proprietary model to generate 150 template pairs, totaling 1200 templates. These templates are minimally different in semantics but vary across a broad range of topics, verb choices, names, and syntactic structures. Examples are presented in Table 7. For instance, a representative pair might include:

- [name] has  $\{x\}$  apples. They get  $\{y\}$  more apples. How many apples does [name] have now? (corresponding to  $x + y$ )
- [name] has  $\{x\}$  apples. They give away  $\{y\}$  apples. How many apples does [name] have now? (corresponding to  $x - y$ )

Each template is instantiated by replacing the [name] placeholder with a randomly selected name from a curated list of 30 English first names, shown below:

James, Emma, William, Olivia, Benjamin, Charlotte, Henry, Amelia, Alexander, Ava, Samuel, Sophia, Jacob, Mia, Daniel, Lily, Michael, Grace, Ethan, Ella, Jack, Chloe, Lucas, Harper, Thomas, Zoe, Matthew, Nora, Nathan, Isla.

The numerical placeholders  $\{x\}$  and  $\{y\}$  are populated with integers  $\leq 50$ , to avoid detokenization issues during model processing.

## B.2 Logit Attribution and Activation Patching Additional Results

**Other Models** We observe a similar *abstract-then-compute* mechanism in other models, including Qwen 2.5 7B and Qwen 2.5 14B. In Qwen 2.5 7B, the abstraction stage occurs around layers 18–20, with the computation stage beginning around layers 22–23. In Qwen 2.5 14B, abstraction takes place around layers 29–32, followed by computation starting at layer 36.

For additional interpretability results using logit lens and activation patching:

- See Figure 12, Figure 13, and Figure 14 for Llama-3 8B on subtraction, multiplication, and division.
- See Figure 15, Figure 16, Figure 17, and Figure 18 for Qwen 2.5 7B on all four operations: addition, subtraction, multiplication, and division.
- See Figure 19, Figure 20, Figure 21, and Figure 22 for Qwen 2.5 14B on the same set of arithmetic tasks.



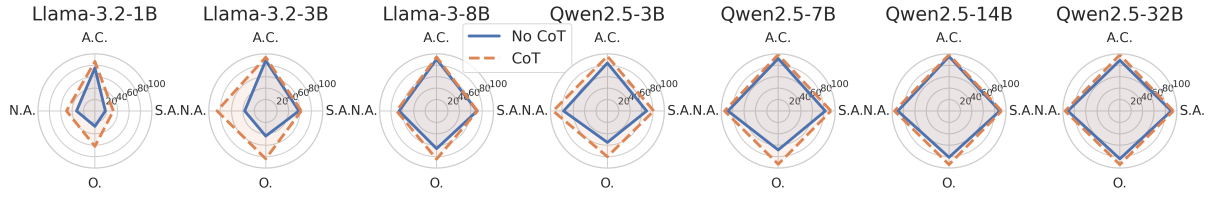


Figure 11: Model zero-shot **with and without CoT** performance on SVAMP. A.C.: Arithmetic Computation; N.A.: Numerical Abstraction; O.: Original; S.A.: Symbolic Abstraction.

Subset	Example Data
(+, -)	(+) [name] owns x stuffed animals. A relative sends them y more stuffed animals. How many stuffed animals does [name] have now? (-) [name] owns x stuffed animals. They give y stuffed animals to a younger sibling. How many stuffed animals does [name] have now?
(+, -)	(+) [name] finds x seashells at the beach. The next day they find y more seashells. How many seashells does [name] have now? (-) [name] finds x seashells at the beach. The tide washes away y seashells. How many seashells does [name] have now?
(+, -)	(+) The storage has x gigabytes free. [name] saves y gigabytes of photos. How much space remains? (-) The storage has x gigabytes free. Cloud storage adds y gigabytes. What is the new capacity?
(×, ÷)	(×) The glacier recedes x inches daily. How much will it shrink after y days? (÷) The glacier retreated x inches over y days. What was the average daily recession?
(×, ÷)	(×) Each server rack uses x kilowatts. What’s the total power for y racks? (÷) The data center used x kilowatts across y racks. What was the average per rack?
(×, ÷)	(×) The spaceship’s shield blocks x radiation units hourly. How much radiation can it block in y hours? (÷) The shield blocked x units over y hours. What was its average protection rate?
Two operations	( $x + y + z$ ) [name] collects x stamps, buys y more, and inherits z. Total stamps? ( $x + y - z$ ) [name] has x stamps, acquires y more, but loses z. How many left? ( $x - y + z$ ) [name] owns x stamps, sells y, but trades for z. How many now? ( $x - y - z$ ) [name] has x stamps, donates y, and ruins z. How many remain?

Table 7: Interpretability dataset examples.

**Two-Operator Dataset** For two operator dataset, we only report results for Qwen 2.5 7B and Qwen 2.5 14B, because Llama-3 8B only achieve 16.5% accuracy on this dataset.

See Figure 27 and Figure 28 for logit attribution results for Qwen 2.5 7B and Qwen 2.5 14B, respectively.

### B.3 Cross-Prompt Patching Additional Results

**Other Models** See Figure 23, Figure 24, and Figure 25 for symbolic abstraction cross-prompt patching results (for single operators: +, -, ×, ÷) on Llama3 8B, Qwen 2.5 7B, and Qwen 2.5 14B, respectively. The results are consistent across models: the likelihood of the target answer peaks at the abstraction stage, while the likelihood of the corrupted answer drops significantly starting from the same stage.

See Figure 26 for numerical abstraction cross-

prompt patching results on Llama3 8B, Qwen 2.5 7B, and Qwen 2.5 14B. We observe consistent trends across all models: the probability of the target answer begins to rise at the onset of the abstraction stage and peaks by its end. Meanwhile, the clean answer probability increases steadily throughout the abstraction stage, reaching a log-probability of 0 at the start of the computation stage.

**Two-Operator Dataset** For the two-operator dataset, we report results only for Qwen 2.5 7B and Qwen 2.5 14B, as Llama-3 8B performs poorly on this setting, achieving only 16.5% accuracy.

See Figure 29 for symbolic abstraction cross-prompt patching results on Qwen 2.5 7B and Qwen 2.5 14B.

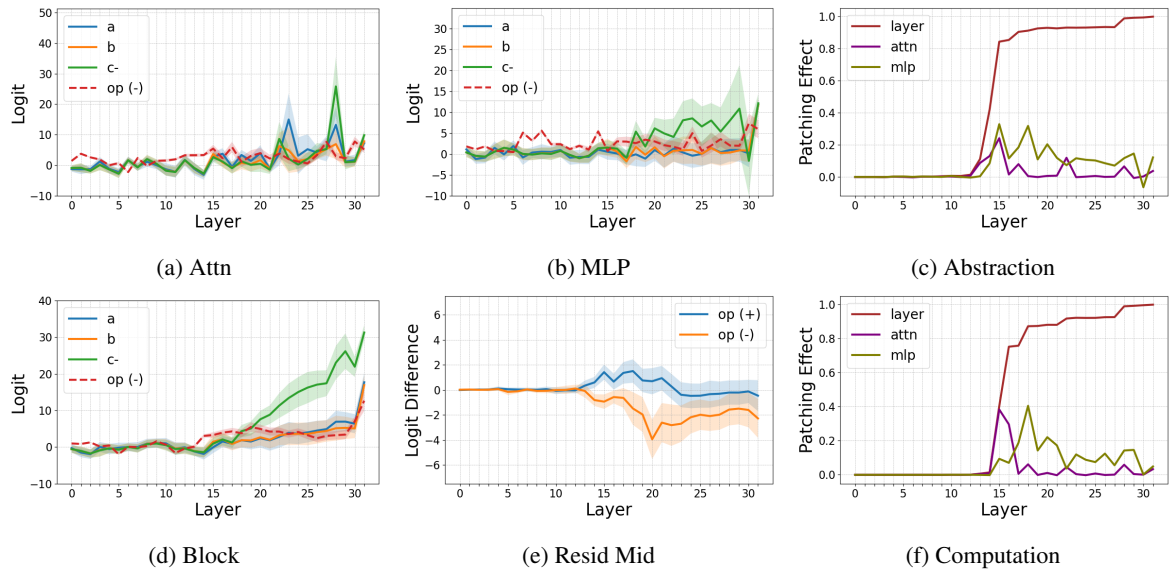


Figure 12: Visualizations of internal computations at last token position in **Llama-3 8B** for **subtraction** math word problems: (a, b, d, e) for logit attribution results, (c, d) activation patching for results.

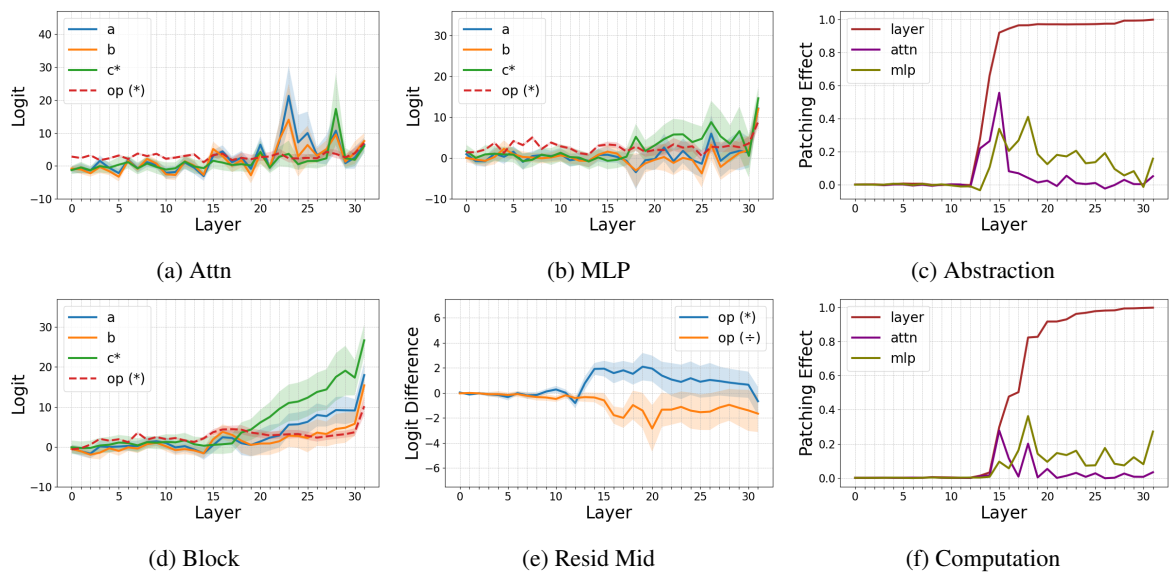


Figure 13: Visualizations of internal computations at last token position in **Llama-3 8B** for **multiplication** math word problems: (a, b, d, e) for logit attribution results, (c, d) activation patching for results.

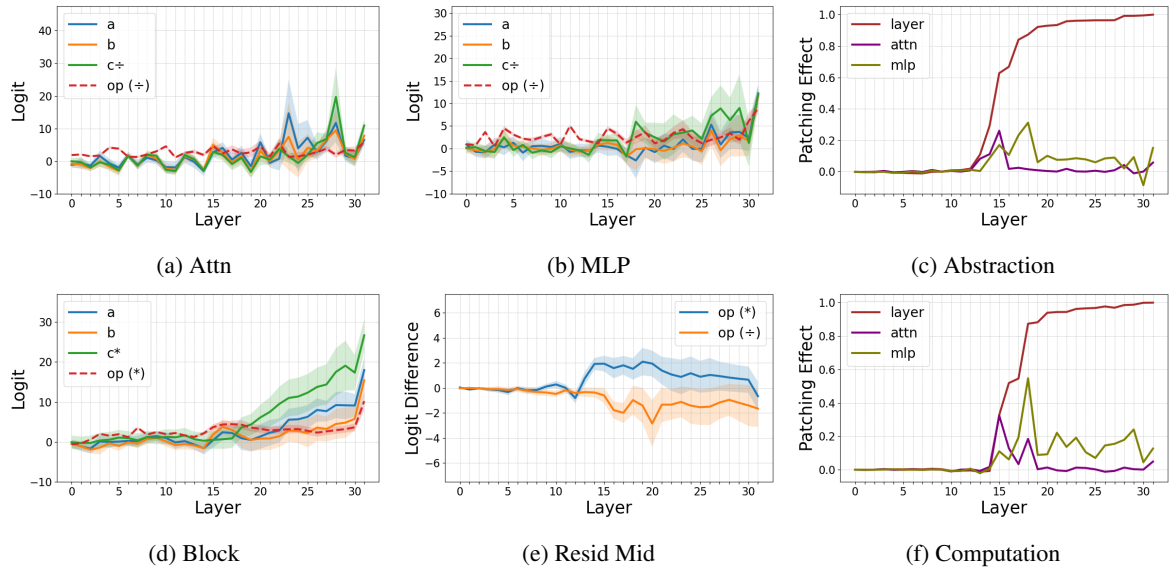


Figure 14: Visualizations of internal computations at last token position in **Llama-3 8B** for **division** math word problems: (a, b, d, e) for logit attribution results, (c, d) activation patching for results.

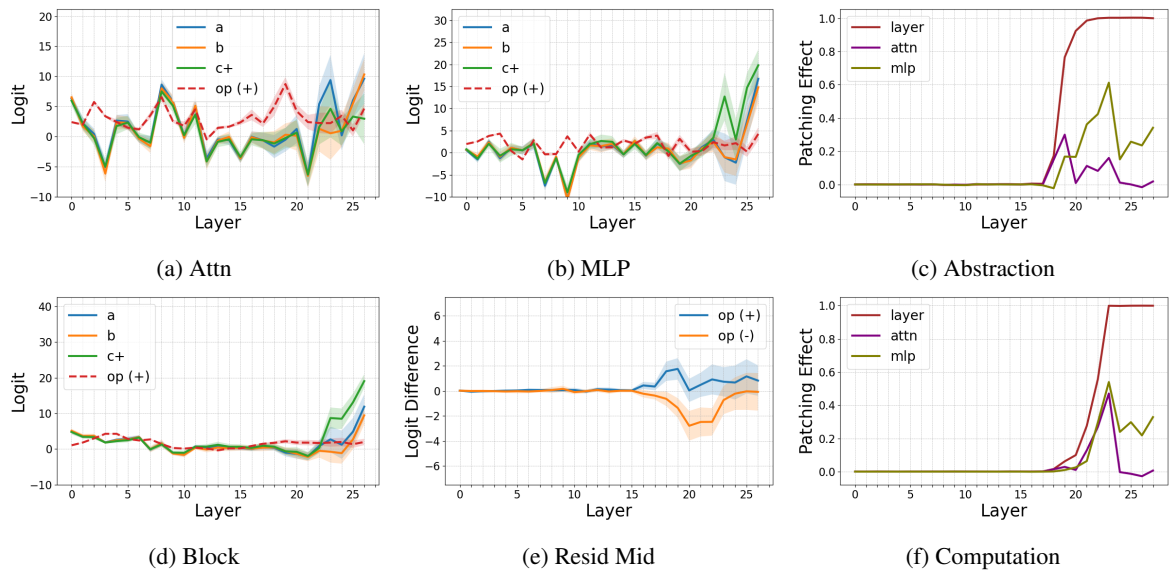


Figure 15: Visualizations of internal computations at last token position in **Qwen 2.5 7B** for **addition** math word problems: (a, b, d, e) for logit attribution results, (c, d) activation patching for results.

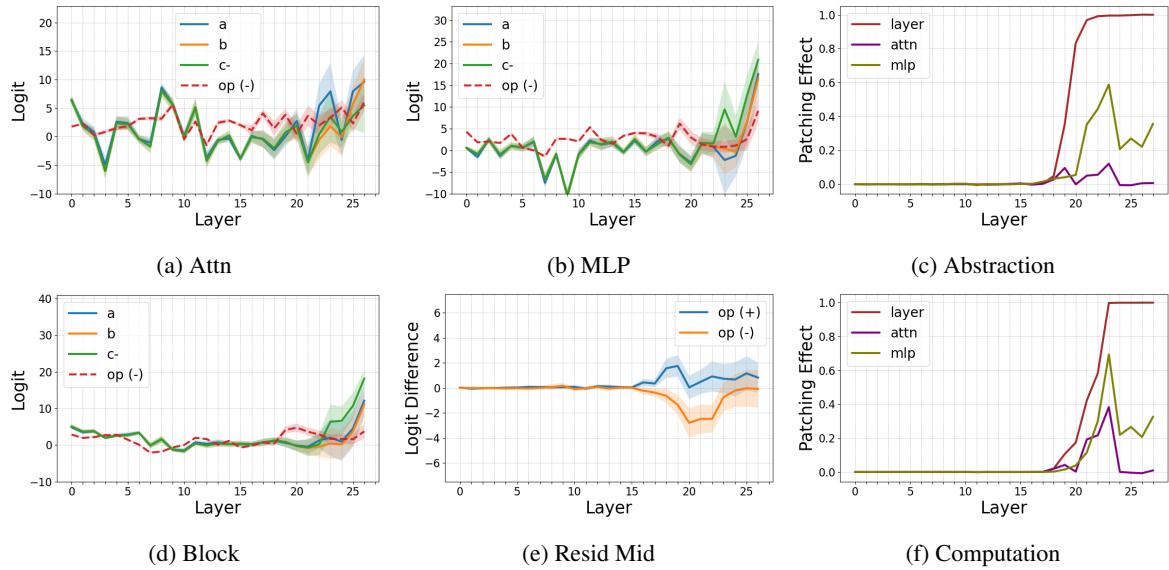


Figure 16: Visualizations of internal computations at last token position in **Qwen 2.5 7B** for **subtraction** math word problems: (a, b, d, e) for logit attribution results, (c, d) activation patching for results.

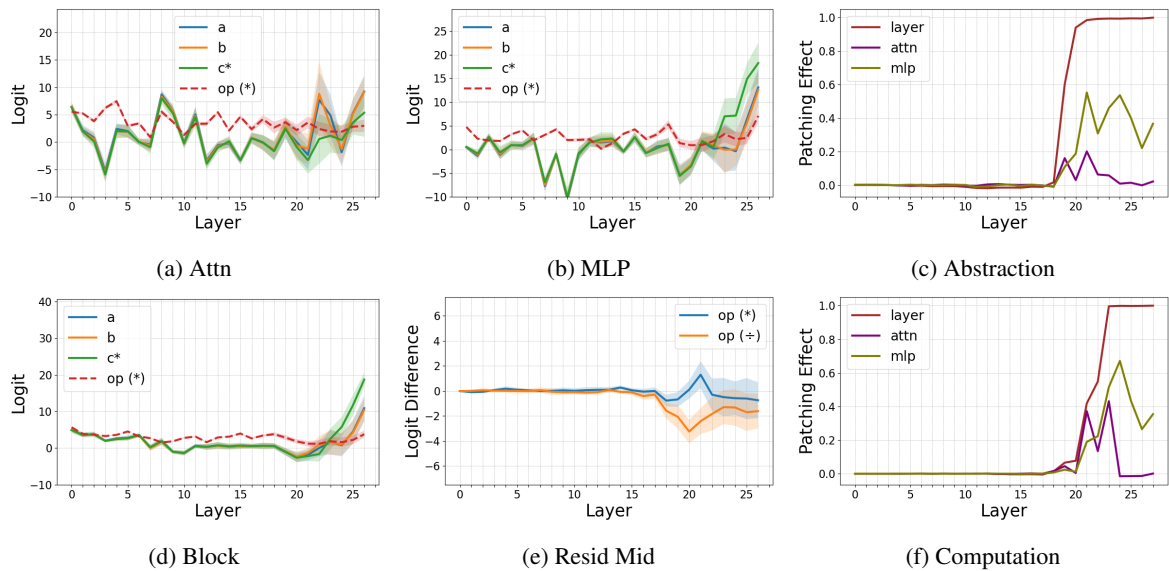


Figure 17: Visualizations of internal computations at last token position in **Qwen 2.5 7B** for **multiplication** math word problems: (a, b, d, e) for logit attribution results, (c, d) activation patching for results.

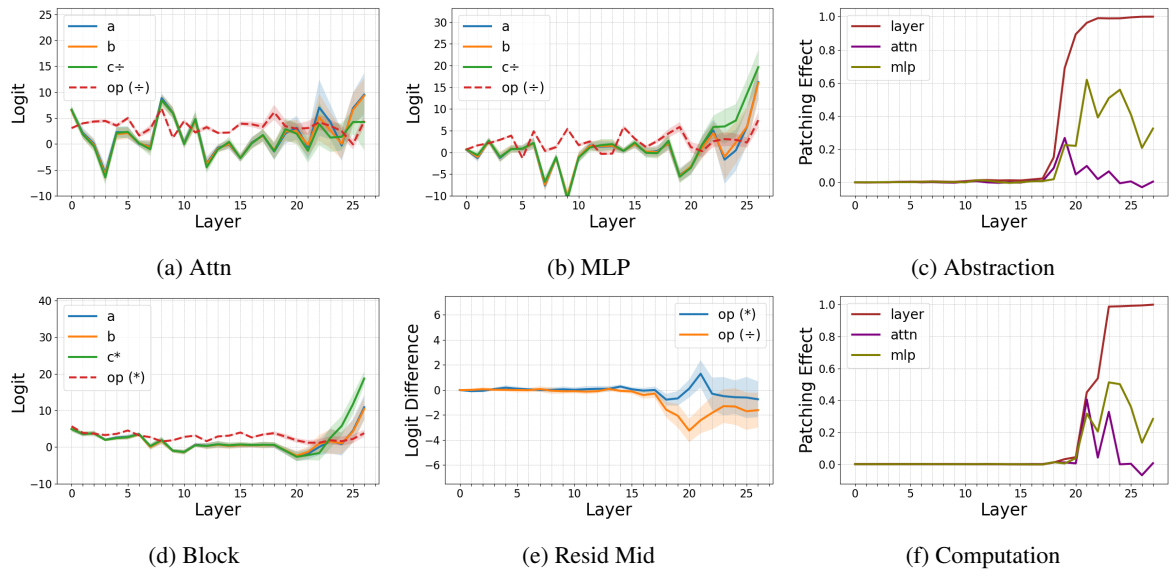


Figure 18: Visualizations of internal computations at last token position in **Qwen 2.5 7B** for **division** math word problems: (a, b, d, e) for logit attribution results, (c, d) activation patching for results.

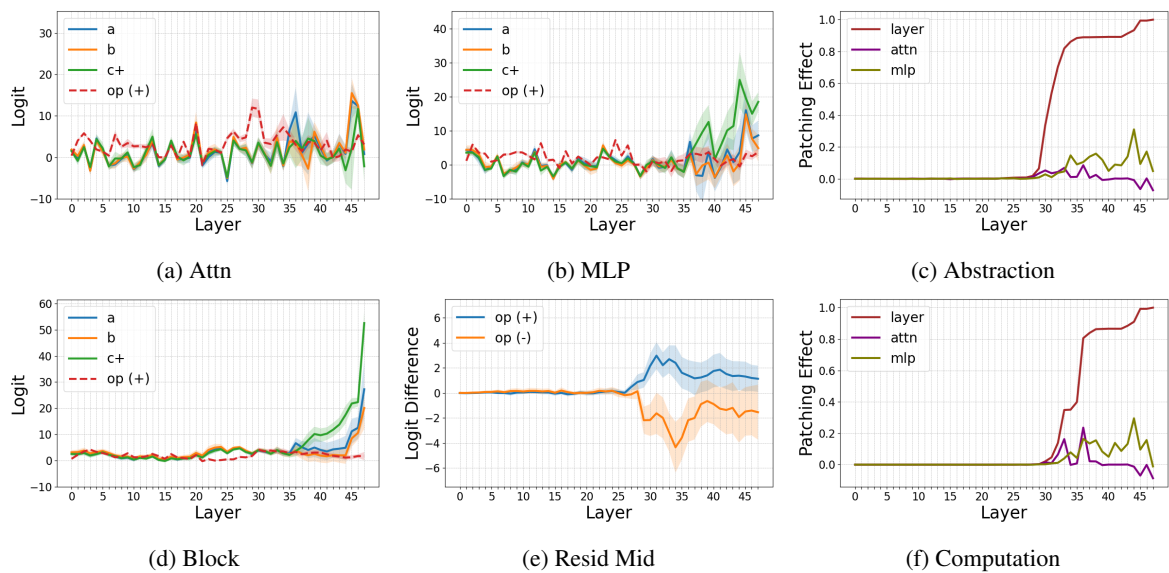


Figure 19: Visualizations of internal computations at last token position in **Qwen 2.5 14B** for **addition** math word problems: (a, b, d, e) for logit attribution results, (c, d) activation patching for results.

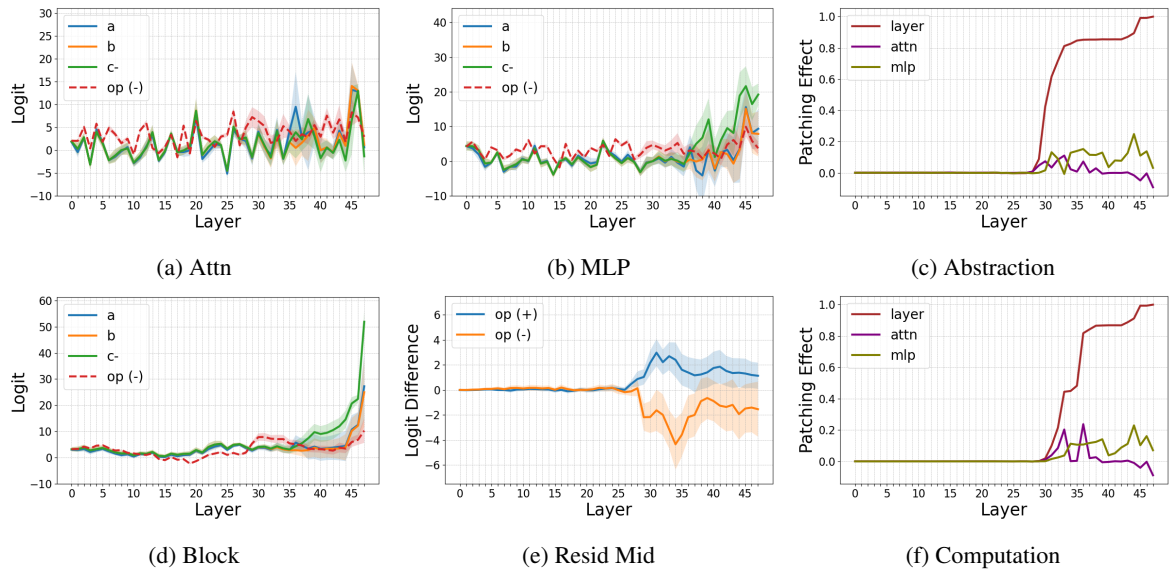


Figure 20: Visualizations of internal computations at last token position in **Qwen 2.5 14B** for **subtraction** math word problems: (a, b, d, e) for logit attribution results, (c, d) activation patching for results.

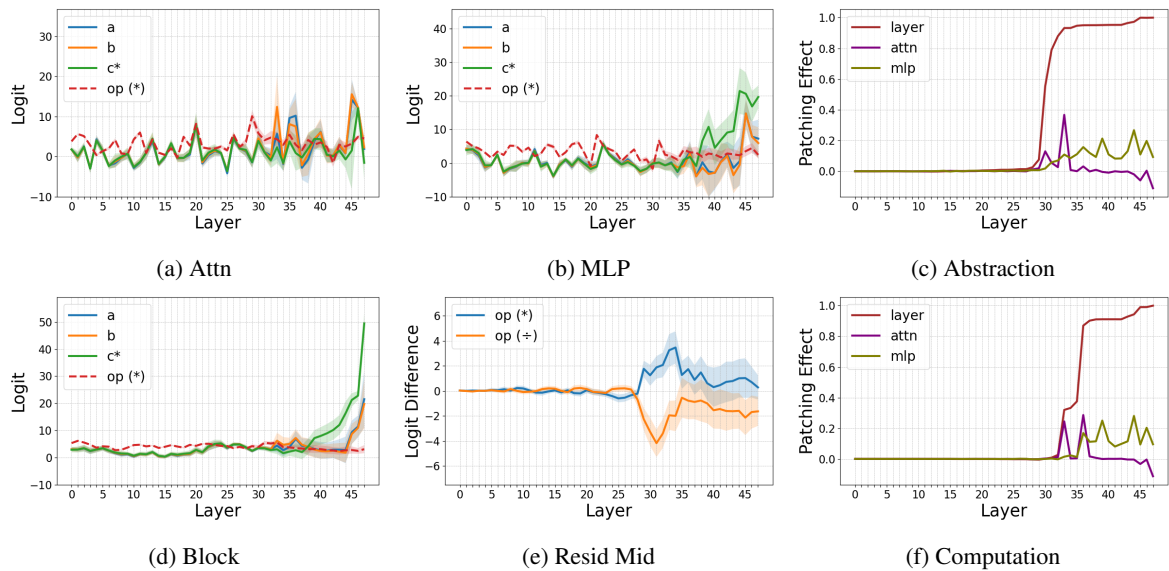


Figure 21: Visualizations of internal computations at last token position in **Qwen 2.5 14B** for **multiplication** math word problems: (a, b, d, e) for logit attribution results, (c, d) activation patching for results.

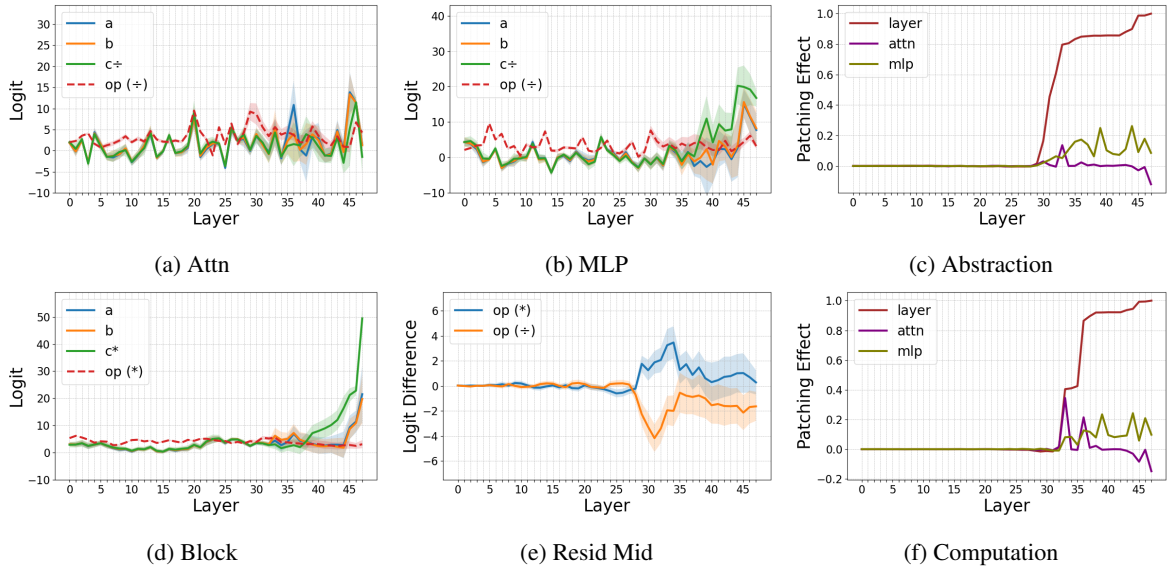


Figure 22: Visualizations of internal computations at last token position in **Qwen 2.5 14B** for **division** math word problems: (a, b, d, e) for logit attribution results, (c, d) activation patching for results. We label the starting layer of abstraction, operand moving and computation in pink, blue and green, respectively.

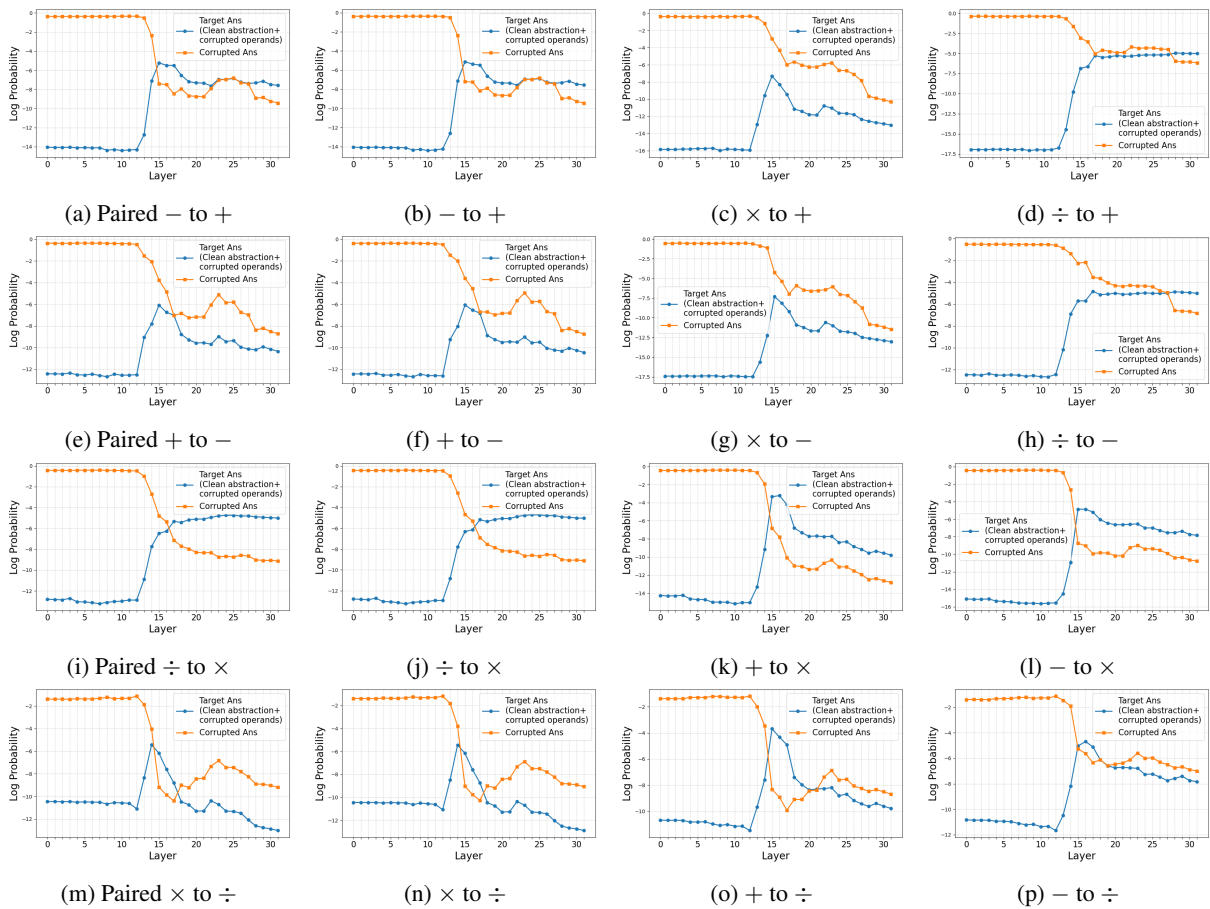


Figure 23: **Llama-3 8B** cross-prompt patching for **symbolic abstraction** results: **First row**: patching symbolic logic to concrete **addition**; **Second row**: patching symbolic logic to concrete **subtraction**; **Third row**: patching symbolic logic to concrete **multiplication**; **Fourth row**: patching symbolic logic to concrete **division**;

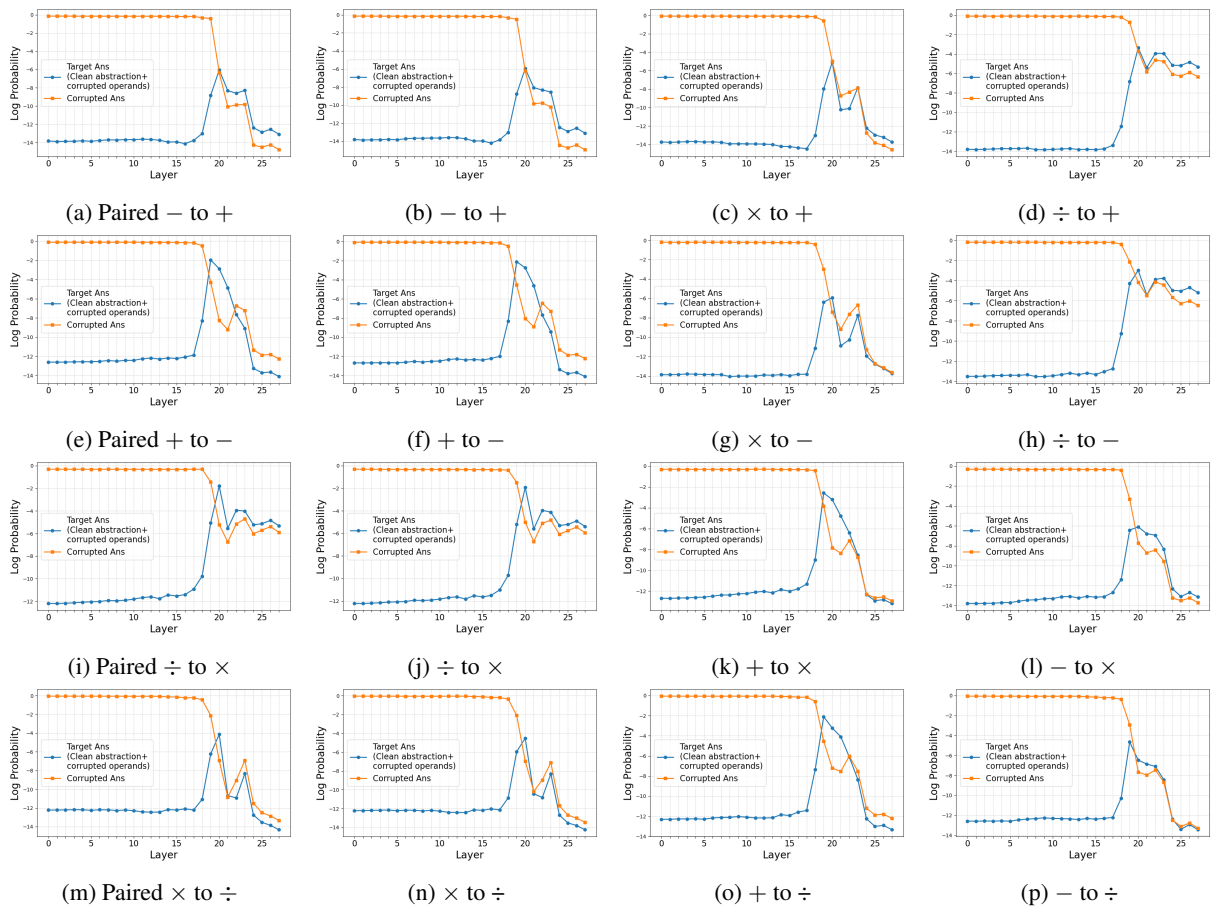


Figure 24: **Qwen-7b** cross-prompt patching for **symbolic abstraction** results: **First row**: patching symbolic logic to concrete **addition**; **Second row**: patching symbolic logic to concrete **subtraction**; **Third row**: patching symbolic logic to concrete **multiplication**; **Fourth row**: patching symbolic logic to concrete **division**;



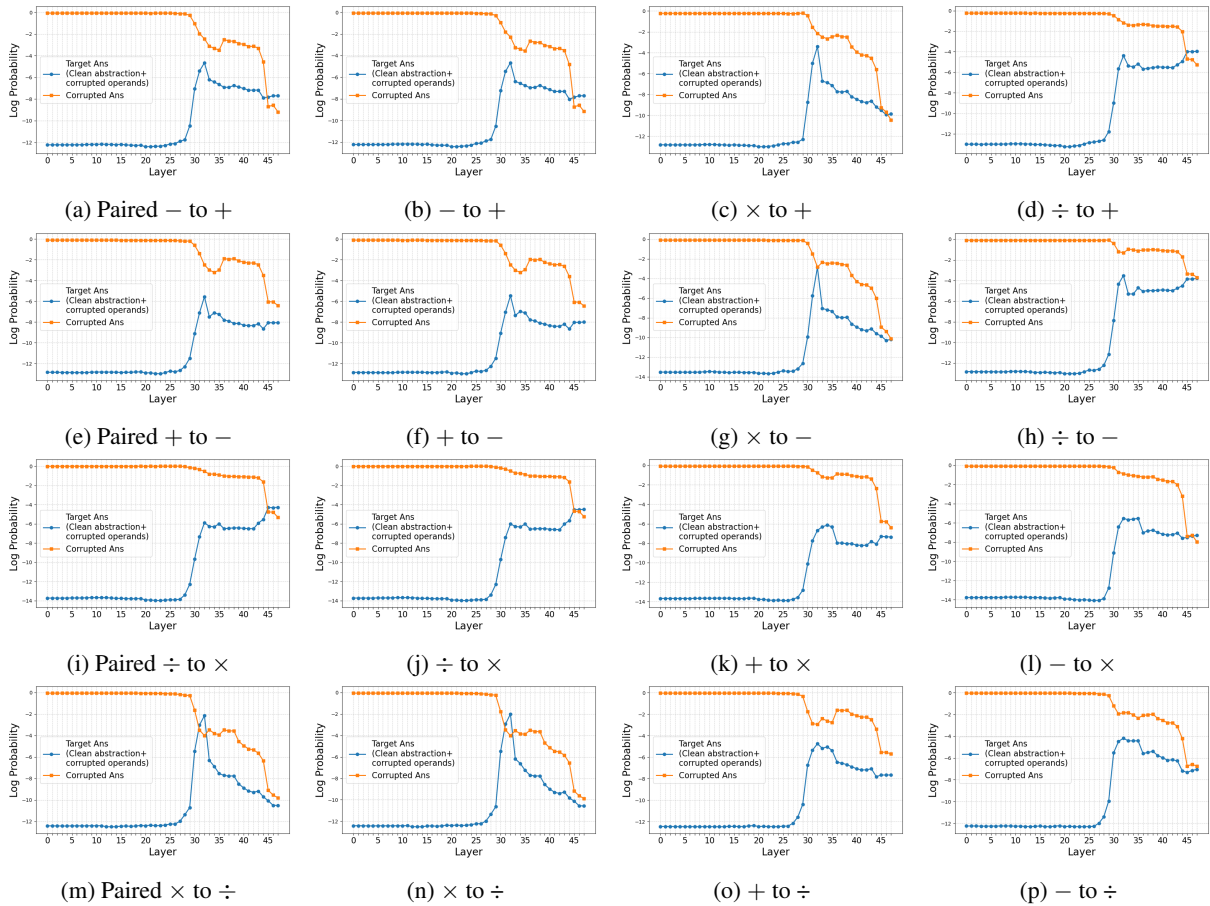


Figure 25: **Qwen-14b** cross-prompt patching for **symbolic abstraction** results: **First row**: patching symbolic logic to concrete **addition**; **Second row**: patching symbolic logic to concrete **subtraction**; **Third row**: patching symbolic logic to concrete **multiplication**; **Fourth row**: patching symbolic logic to concrete **division**;

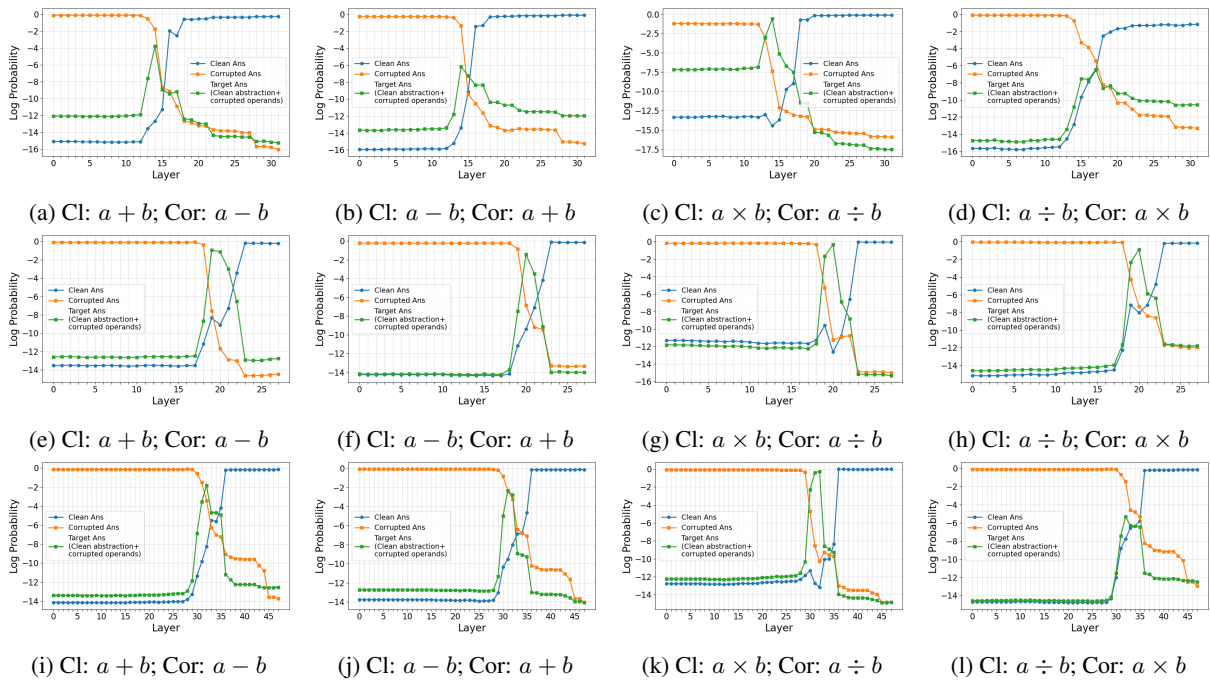


Figure 26: Cross-prompt patching results for **numerical abstraction**. **First row**: results for **Llama-3 8B** with corresponding clean and corrupted run. **Second row**: results for **Qwen2.5 7B** with corresponding clean and corrupted run. **Third row**: results for **Qwen2.5 14B** with corresponding clean and corrupted run.

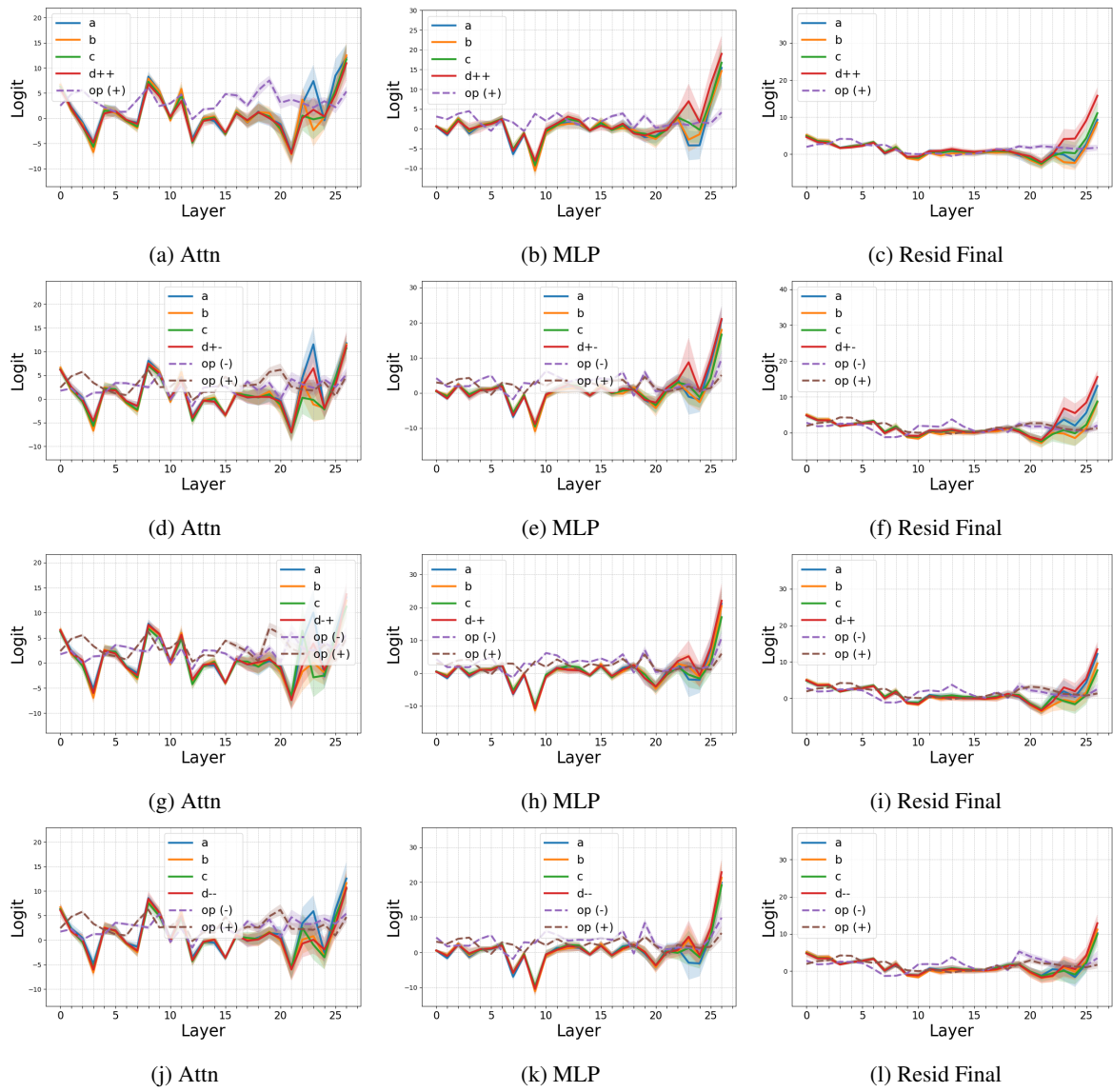


Figure 27: Visualizations of internal computations at last token position in **Qwen 2.5 7B** for **two-operation** math word problems. **First row:** for  $a + b + c$ . **Second row:** for  $a + b - c$ . **Third row:** for  $a - b + c$ . **Fourth row:** for  $a - b - c$ .

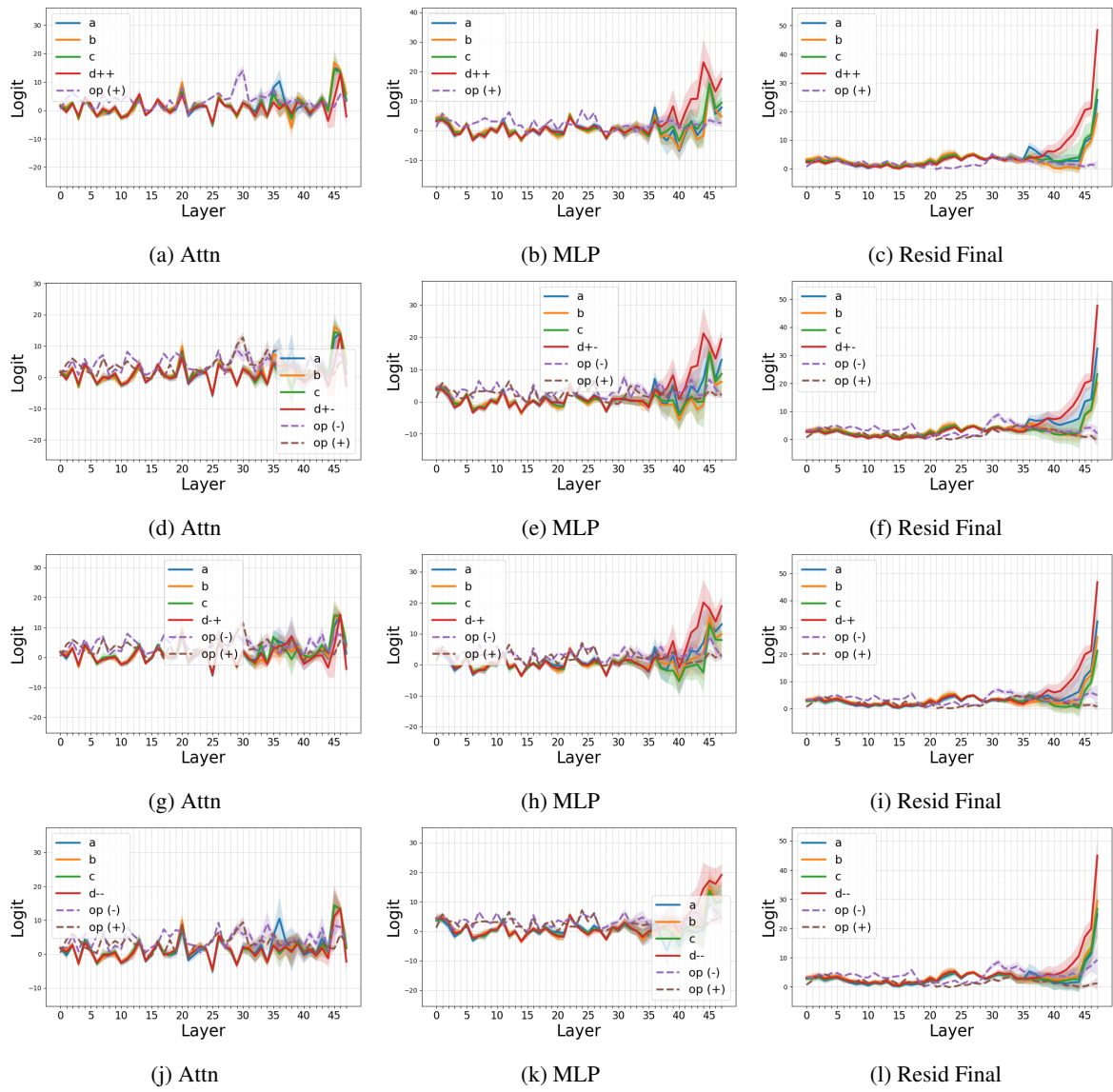


Figure 28: Visualizations of internal computations at last token position in **Qwen 2.5 14B** for **two-operation** math word problems. **First row:** for  $a + b + c$ . **Second row:** for  $a + b - c$ . **Third row:** for  $a - b + c$ . **Fourth row:** for  $a - b - c$ .

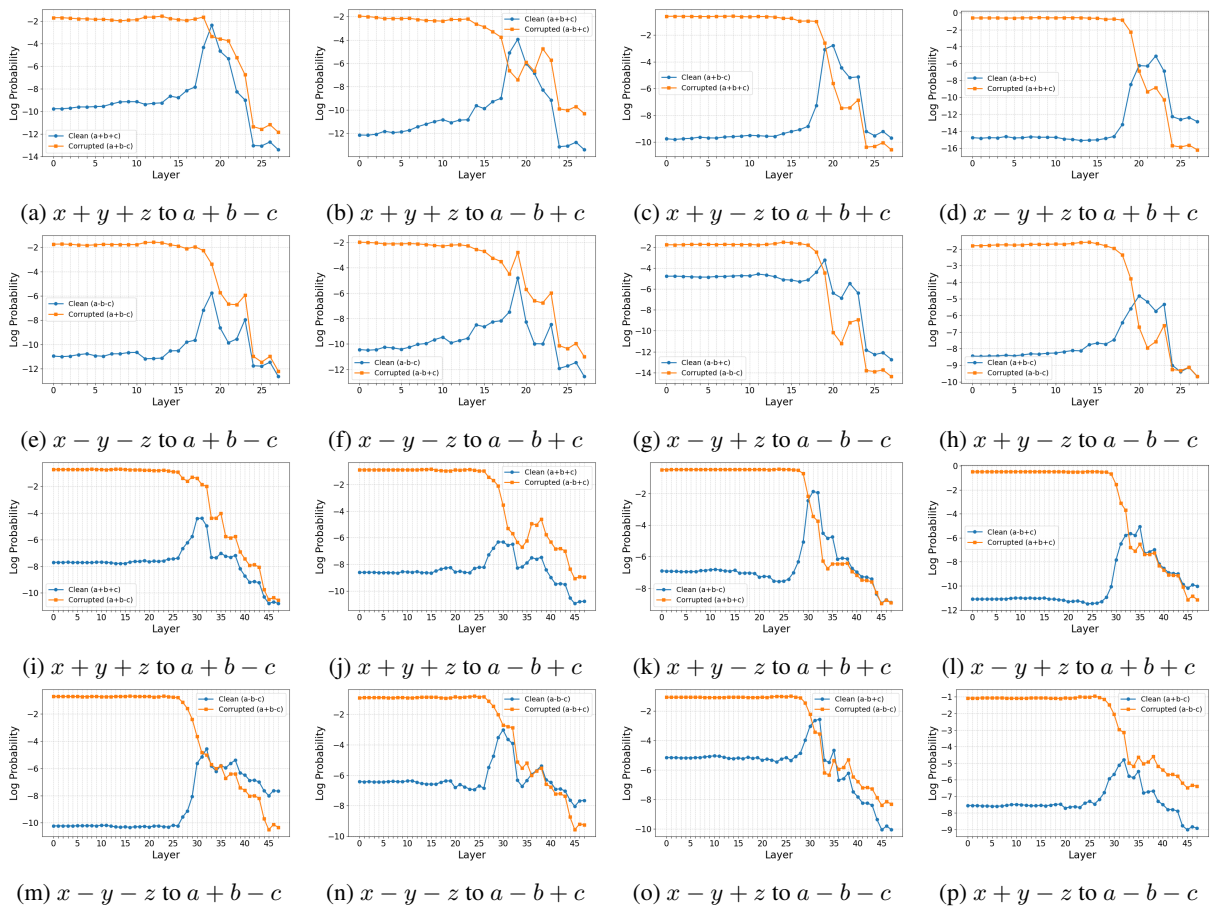


Figure 29: **Two-operation** cross-prompt patching for **symbolic abstraction** results: **First row & Second row:** patching symbolic logic to concrete problems for Qwen 2.5 7B. **Third row & Fourth row:** patching symbolic logic to concrete problems for Qwen 2.5 14B.