# Think-to-Talk or Talk-to-Think?
# When LLMs Come Up with an Answer in Multi-Step Reasoning

**Keito Kudo[1,2], Yoichi Aoki[1,2], Tatsuki Kuribayashi[3], Shusaku Sone[1]**
**Masaya Taniguchi[2], Ana Brassard,[2,1] Keisuke Sakaguchi[1,2], Kentaro Inui[3,1,2]**
[1]Tohoku University, [2]RIKEN, [3]MBZUAI
{keito.kudo.q4, youichi.aoki.p2, sone.shusaku.r8}@dc.tohoku.ac.jp,
{tatsuki.kuribayashi, kentaro.inui}@mbzuai.ac.ae,
keisuke.sakaguchi@tohoku.ac.jp, {masaya.taniguchi, ana.brassard} @riken.jp

## Abstract

This study investigates the internal reasoning mechanism of language models during symbolic multi-step reasoning, motivated by the question of whether chain-of-thought (CoT) outputs are faithful to the model's internals. Specifically, we inspect *when* they internally determine their answers, particularly before or after CoT begins, to determine whether models follow a post-hoc "think-to-talk" mode or a step-by-step "talk-to-think" mode of explanation. Through causal probing experiments in controlled arithmetic reasoning tasks, we found systematic internal reasoning patterns across models; for example, simple subproblems are solved before CoT begins, and more complicated multi-hop calculations are performed during CoT.[1]

## 1 Introduction

A lengthy explanation may be produced in two modes: as a post-hoc explanation to a predetermined conclusion (*think-to-talk*) or by process of reaching a conclusion while explaining (*talk-to-think*). An analogy applies to large language models (LLMs) using chain-of-thought (CoT; Wu et al., 2023)-style reasoning: *is the output a post-hoc explanation, or does it reflect step-by-step solving?*

To answer this question, we applied causal probing to model internals at each layer during each timestep to determine *when* the answer is reached in CoT-style reasoning. This work joins a line of research concerned with a mechanistic interpretation of LLMs (Conneau et al., 2018; Tenney et al., 2019; Niven and Kao, 2019; nostalgebraist, 2020; Li et al., 2023; Heinzerling and Inui, 2024; Lieberum et al., 2024; Yang et al., 2024c, *i.a.*) and the practical necessity of controlling their behavior (Zhao et al., 2019; Ganguli et al., 2023; Yang et al., 2024b, *i.a.*). Specifically, we extend the target of existing analysis from simple arithmetics, such as by Stolfo et al.
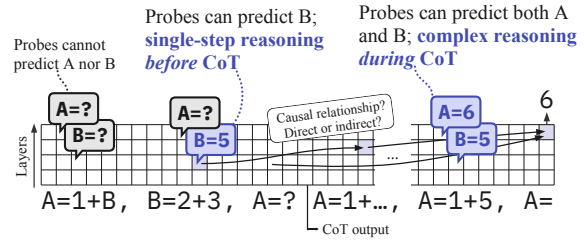


Figure 1: Probes predict each variable's value from the LLMs' hidden states at each timestep×layer; a high accuracy indicates where the model completes the corresponding calculation (§ 3). Causal interventions (§ 4.1, § 4.2) revealed a somewhat indirect relationship between sub- and final answers.

(2023), to *multi-step* reasoning outputs, making it possible to differentiate between simple and more complex reasoning steps, as well as finding systematic evidence of reasoning being conducted both before and *during* CoT.

In our experiments, we prepared a controlled testbed of symbolic arithmetic reasoning tasks and observed whether trained classifiers, given a set of hidden representations, could accurately predict the values of all variables in the test instance. By comparing accuracies across each timestep, we can observe at which point the model's internals start being informative to the probes, illustrating the model's internal reasoning flow. We repeated this analysis across ten LLMs.

Overall, we found that simple single-step reasoning (e.g., A=2+3,A=?) is performed before CoT, and more complex reasoning (e.g., A=2+3,B=A+3,B=?) follows after CoT (§ 3; Fig 1). Furthermore, causal interventions revealed that such predetermined (sub-)answers impacted the model's final answer (§ 4.1), but their causal relationship was somewhat indirect (§ 4.2), elucidating a mixed but systematic internal reasoning pattern between *think-to-talk* and *talk-to-think* modes.

---

[1]The code/data will be made public upon acceptance.

| Level | INPUT | OUTPUT | #Step | #Stack | #Dist. |
|---|---|---|---|---|---|
| 1 | $\underline{A = 1 + B}_{-3}, \underline{B = 2}_{-2}; \underline{A =?}_{-1}$ | $\underline{A = 1 + B}_0, \underline{B = 2}_1, \underline{A = 1 + B}_2, \underline{A = 1 + 2}_3,$ $\underline{A = 3}_4$ | 1 | 1 | 0 |
| 2 | $\underline{A = 2 + 3}_{-3}, \underline{B = 1 + A}_{-2}; \underline{B =?}_{-1}$ | $\underline{B = 1 + A}_0, \underline{A = 2 + 3}_1, \underline{A = 5}_2, \underline{B = 1 + A}_3,$ $\underline{B = 1 + 5}_4, \underline{B = 6}_5$ | 2 | 0 | 0 |
| 3 | $\underline{A = 1 + B}_{-3}, \underline{B = 2 + 3}_{-2}; \underline{A =?}_{-1}$ | $\underline{A = 1 + B}_0, \underline{B = 2 + 3}_1, \underline{B = 5}_2, \underline{A = 1 + B}_3,$ $\underline{A = 1 + 5}_4, \underline{A = 6}_5$ | 2 | 1 | 0 |
| 4 | $\underline{A = 1 + B}_{-4}, \underline{B = 2 + 3}_{-3}, \underline{C = 4 + 5}_{-2}; \underline{A =?}_{-1}$ | $\underline{A = 1 + B}_0, \underline{B = 2 + 3}_1, \underline{B = 5}_2, \underline{A = 1 + B}_3,$ $\underline{A = 1 + 5}_4, \underline{A = 6}_5$ | 2 | 1 | 1 |
| 5 | $\underline{A = 1 + B}_{-4}, \underline{B = 2 + C}_{-3}, \underline{C = 1 + 2}_{-2}; \underline{A =?}_{-1}$ | $\underline{A = 1 + B}_0, \underline{B = 2 + C}_1, \underline{C = 1 + 2}_2, \underline{C = 3}_3,$ $\underline{B = 2 + C}_4, \underline{B = 2 + 3}_5, \underline{B = 5}_6, \underline{A = 1 + B}_7,$ $\underline{A = 1 + 5}_8, \underline{A = 6}_9$ | 3 | 2 | 0 |

Table 1: Examples of arithmetic reasoning tasks employed in our experiments. We set 5 levels with different complexities, e.g., #Step number. Specifically, #Step indicates the number of required operations to reach the final answer. #Stack indicates how many variables' values are not immediately determined in their first appearing equation. #Dist. is the number of unnecessary distractor equations. The equation position (e.g., $-\mathbf{3}$) is noted in each equation's lower right corner, which is referred to by the $\mathrm{RSP}_{\mathrm{eq}}$ score and $\mathrm{eq}(\cdot)$ function in Eq.2.

## 2 Experimental settings

### 2.1 Evaluation task

**Dataset:** We adopt five settings with different complexity levels, e.g., the number of required reasoning steps to reach the answer (see Table 1); these data are borrowed from Kudo et al. (2023). Each problem consists of equations with primitive operations, such as reference (A=B), assignment (B=1), and arithmetic operation (C=1+3), finally asking for the value of a particular variable $v_n$. We create 12,000 instances in each level (10,000 for training probing classifiers and 2,000 for testing their accuracy; see § 2.2) by changing the variable names and numbers appearing in the equations (see Appendix A for details).

**Model inference:** In each task level, given the INPUT $x$, models must generate an OUTPUT consisting of intermediate steps $z$ and a final answer $y$, e.g., A=3. The task is formatted in a few-shot setting that provides three examples with the exact same levelAs a sanity check, we confirmed that our used models could follow the OUTPUT formant and solve the tasks with nearly 100% accuracy.

**Research focus:** As stated in § 1, our question is when models solved the (sub)problems in the CoT-style reasoning. A key to answering this question will be *where it is possible to extract the final answer (or necessary information for it) from the model's internal representations*—we observe reasoning patterns across task levels and models (§ 3).

### 2.2 Probing

For each task level, we train a linear probe (Alain and Bengio, 2017) for each combination of token position $t$, layer number $l$, and $i$-th variable $v_i$ in a problem to predict the number finally assigned to $v_i$. Specifically, given a model's hidden state $\boldsymbol{h}_{t,l} \in \mathbb{R}^d$, the probing classifier $f_{t,l,v_i}(\cdot) : \mathbb{R}^d \to \{0, \ldots, 9\}$ predicts the correct value of $v_i$ (Appendix B). If a probe could yield a good accuracy, this suggests that the computation of $v_i$ value is already done at the corresponding position (position $t$ and layer number $l$). Figure 2 shows an example of our probing results at level 3, where, for example, the value of B can be extracted in INPUT and thus already computed before CoT begins. § 4 analyzes the causal relationship between the identified reasoning patterns and the final answer, supporting the validity of our linear probing approach.

### 2.3 Evaluation metrics

The probing results are aggregated as follows:

$$\mathrm{RSP}_{\mathrm{tok}}(v_i) = \min\{t \mid \max_l \mathrm{acc}(t, l, v_i) > \tau\} , \quad (1)$$

$$\mathrm{RSP}_{\mathrm{eq}}(v_i) = \mathrm{eq}\left(\mathrm{RSP}_{\mathrm{tok}}(v_i)\right) , \quad (2)$$

where $\mathrm{acc}(t, l, v_i) \in [0, 1]$ indicates the accuracy of probing classifier $f_{t,l,v_i}$. The $\mathrm{RSP}_{\mathrm{tok}}$ and $\mathrm{RSP}_{\mathrm{eq}}$ (relative spike position; RSP) indicate when is the first time probing classifier achieved a high accuracy above $\tau$ (= 0.90 in our study), and position is denoted as the token or equation index, respectively. $\mathrm{eq}(\cdot) : \mathbb{Z} \to \mathbb{Z}$ in Eq. 2 maps a token index to the equation index it belongs to (see Table 1). The token/equation indices are denoted as a relative position to the start of OUTPUT ($t = \mathrm{eq}(t) = 0$) where CoT begins, i.e., $t \in \mathbb{Z}$ takes a negative value in an INPUT part. The lower bounds of RSP scores (at which position the value of a variable
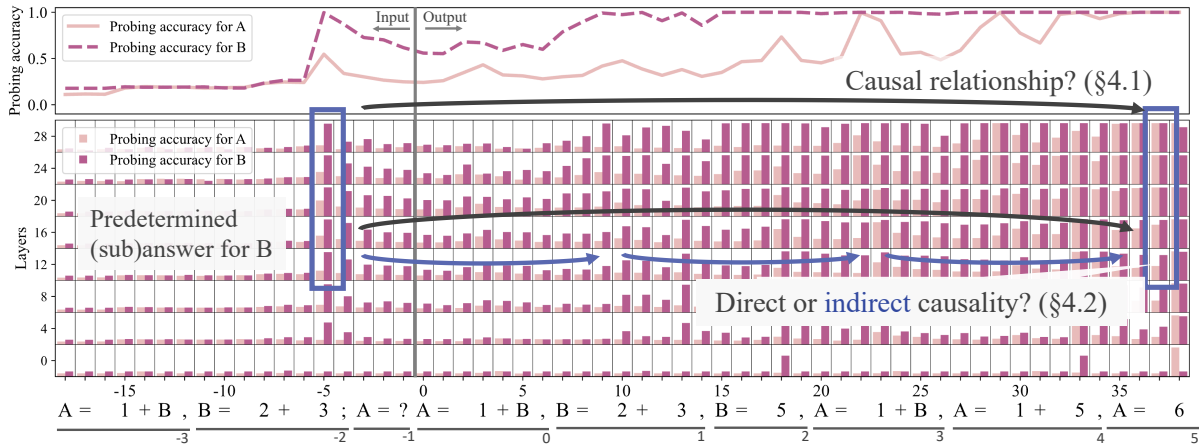
Figure 2: Probing results for Qwen2.5-7B at the task level 3. The input sequence below the graphs is one example (the results are averaged over the test set). The upper part indicates the maximum probing accuracy achieved at each token position $t$. The bottom part shows the probing accuracies in each token $t$, layer $l$, and variable $v_i$. § 4 analyzes their causal relationship with the final answer.

is uniquely determined first) can be determined by a left-to-right greedy solver and reported as a gold baseline in our experiments. In the example of Figure 2, for example, $\mathrm{RSP}_{\mathrm{tok}}(\mathrm{B})$ is a negative value, indicating that B value was computed before CoT begins. For reference, we additionally report three types of accuracy: $\mathrm{Acc}_{<\mathrm{CoT}}$, $\mathrm{Acc}_{>\mathrm{CoT}}$, and task accuracy. $\mathrm{Acc}_{<\mathrm{CoT}}$ and $\mathrm{Acc}_{>\mathrm{CoT}}$ indicate the maximum probing accuracy achieved before and after CoT begins. Task accuracy indicates the exact match accuracy of OUTPUT (i.e., $z$ and $y$).

## 3 Experimental results

**Across task levels:** We first analyze Qwen2.5-7B (Qwen Team, 2024) across the five task levels. Table 2 shows RSP scores as well as how many primitive operations are required to compute the value of $v_i$ (#steps) and their low bounds of $\mathrm{RSP}_{\mathrm{eq}}$ (eq*). The results show that variables with fewer required reasoning steps (#steps≤1) achieved the lower bound of $\mathrm{RSP}_{\mathrm{eq}}$ (bold scores in Table 2), and vice versa. This clear pattern is also reflected in the contrast between $\mathrm{Acc}_{<\mathrm{CoT}}$ and $\mathrm{Acc}_{>\mathrm{CoT}}$ (larger contrast is shown in variables with many required steps). Furthermore, these patterns have emerged robustly across five levels. That is, the model had selectively solved single-step subproblems before CoT began and resolved more complicated multi-step computations during CoT. Appendix D shows the fact that consistent results can be obtained independent of in-context example selections and some sanity checks for our approach.

| | Variable | | RSP (↓) | | Acc (↑) | | |
|---|---|---|---|---|---|---|---|
| Level | name in Table 1 | #steps | eq* | eq | < CoT | > CoT | Task |
| 1 | A | 1 | −2 | **−2** | 99.6 | 100 | 100 |
| | B | 0 | −2 | **−2** | 100 | 100 | |
| 2 | A | 1 | −3 | **−3** | 100 | 100 | 100 |
| | B | 2 | −2 | 3 | 52.1 | 100 | |
| 3 | A | 2 | −2 | 3 | 54.7 | 100 | 100 |
| | B | 1 | −2 | **−2** | 99.7 | 100 | |
| 4 | A | 2 | −3 | 4 | 46.0 | 100 | |
| | B | 1 | −3 | **−3** | 99.2 | 100 | 100 |
| | C | 1 | −2 | **−2** | 99.1 | 34.5 | |
| 5 | A | 3 | −2 | 7 | 25.1 | 100 | |
| | B | 2 | −2 | 5 | 49.0 | 100 | 100 |
| | C | 1 | −2 | **−2** | 99.8 | 100 | |

Table 2: The results of Qwen2.5-7B on five levels. The eq* and eq indicate the lower bound and model's $\mathrm{RSP}_{\mathrm{eq}}$ scores. The < CoT, > CoT, and task accuracies are introduced in § 2.3.

**Across models:** We further analyzed ten models listed in Table 3 on the level 3 task. Almost the same results as Qwen2.5-7B are generally obtained across various models, enhancing the generality of our obtained findings. For more precise comparisons, we also computed the $\mathrm{RSP}_{\mathrm{tok}}$ scores, and larger models could induce the intermediate answer slightly earlier than smaller ones, e.g., lower $\mathrm{RSP}_{\mathrm{tok}}$ for Yi1.5-34B than Yi1.5-9B.

## 4 Causal interventions

We investigate the causal dependencies between the predetermined answers identified by our probing (§ 3) and the models' final answer through causal intervention analysis. We use Qwen2.5-7B and

| Model | variable | RSP (↓) | | Acc (↑) | | |
| | | eq | tok | < CoT | > CoT | Task |
|---|---|---|---|---|---|---|
| Qwen2.5 (7B) | A | 3 | 22 | 54.7 | 100 | 100 |
| (Qwen Team, 2024) | B | −2 | −5 | 99.7 | 100 | |
| Qwen2.5 (14B) | A | 3 | 22 | 72.5 | 100 | 100 |
| (Qwen Team, 2024) | B | −2 | −5 | 100 | 100 | |
| Qwen2.5 (32B) | A | 3 | 21 | 69.7 | 100 | 100 |
| (Qwen Team, 2024) | B | −2 | −5 | 100 | 100 | |
| Qwen2.5-Math (7B) | A | 3 | 22 | 63.4 | 100 | 100 |
| (Yang et al., 2024a) | B | −2 | −5 | 100 | 100 | |
| Yi1.5 (9B) | A | 4 | 32 | 50.0 | 100 | 100 |
| (Young et al., 2024) | B | −2 | −5 | 96.5 | 100 | |
| Yi1.5 (34B) | A | 3 | 25 | 57.7 | 100 | 100 |
| (Young et al., 2024) | B | −2 | −5 | 100 | 100 | |
| Llama3.1 (8B) | A | 3 | 22 | 63.4 | 100 | 100 |
| (Dubey et al., 2024) | B | −2 | −5 | 92.2 | 100 | |
| Llama3.2 (3B) | A | 5 | 36 | 47.9 | 98.0 | 97.6 |
| (Dubey et al., 2024) | B | 2 | 10 | 86.5 | 99.5 | |
| Mistral-Nemo (12B) | A | 5 | 35 | 39.6 | 100 | 99.6 |
| (Mistral AI Team, 2024) | B | −2 | −5 | 91.4 | 100 | |

Table 3: Results for various models on the task level 3. The "tok" column shows the $\text{RSP}_{\text{tok}}$ score, and the other columns are same as Table 2. RSP scores that are the same as their lower bounds are bolded.

task level 3 format in the following analyses.

**Activation patching:** We employ an activation patching (Vig et al., 2020; Meng et al., 2022; Zhang and Nanda, 2024), which is one widely adopted technique in mechanistic interpretability research. Let us concisely introduce the method in our context: to inspect the causal relationship between a specific hidden state $h_{t,l}$ and a final answer $y$, the state is replaced with another one $\tilde{h}_{t,l}$ obtained from the same model but with a different input context $\tilde{x}$. Then, if the model's subsequent answer changed to be consistent with $\tilde{x}$ rather than $x$, one can confirm their causal dependence. More concretely, in Figure 2, $h_{-5,12}$ seems to have the information of B value. Its causal relationship with final answer A=6 will be confirmed if the answer changes to, e.g., A=7 by replacing $h_{-5,12}$ with other one $\tilde{h}_{-5,12}$ with a modified context of $\tilde{x}$ (A=1+B,B=2+4...) instead of the original one $x$ (A=1+B,B= 2+3...).

## 4.1 Causality between prior and final answer

We compare the model's final answer $y$ (A=?) when applying the intervention of B value to the hidden states associated with high probing accuracy for B ($H^+$) and those with low accuracy ($H^-$); see Appendix E for details. We observed that the intervention to $H^+$ always altered the final answer A, while the intervention to $H^-$ never changed the

answer, corroborating the general causal dependencies between the predetermined answers B elucidated by our probes (§ 3) and the models' final answer A. In this analysis, the model is generating the OUTPUT part as a free-style generation (same as § 3). Appendix E shows more detailed results.

## 4.2 What happens when predetermined answer conflicts with reasoning chain

In some practical scenarios, e.g., model ensembles, a reasoning chain from another information source can be inputted into the model, which potentially conflicts with its own predetermined answer. We additionally analyze how strongly the model sticks to its own predetermined answers against such a conflict. To test such cases, we apply the activation patching *only* to the B value in the IN-PUT part (e.g., B=2+3... → B=2+4...) by intervening the hidden states $H^+$ same as 4.1, while teacher-forcing the original OUTPUT (B=2+3...) as a subsequent context, which contradicts with the intervened INPUT. Interestingly, the generated final answer was *always* consistent with the original OUTPUT (B=2+3...) rather than the intervened INPUT (B=2+4...) over all the test instances.

To sum up, the model's prior decisions indeed biased the final answer (§ 3) but can be easily overwritten by new information in the context between the input and the final answer parts. In this sense, the model's final answer does not directly refer to its own internal prior answers (like the red arrow in table 2) and thus is not always faithful to them. This can be seen as two-stage reasoning—they first selectively resolve some simple subproblems and internally combine/refine/decline them step-by-step through the CoT process.

## 5 Conclusions

We conducted causal probing analyses of when (sub-)answers are determined in the CoT process, using synthetic arithmetic problems as a controlled testbed. We found systematic evidence of internal reasoning patterns—simple subproblems are solved before the CoT begins, and more complicated multi-step calculations are performed during CoT. Our results extend our understanding of LLMs' symbolic reasoning process during multi-step reasoning, providing evidence for both a post-hoc "think-to-talk" mode and step-by-step "talk-to-think" mode of explanation.

## Limitations

**Comprehensiveness of experimental settings**
Some experiments were conducted with a limited scope; for example, the experiments with various models in § 3 are conducted only on the level 3 task. Additionally, causal interventions (§ 4) are performed only with Qwen2.5-7B. Conducting our experiment with more models and tasks will further enhance the generalizability of the results, although some concerns (e.g., robustness to different in-context examples) were addressed in the Appendix.

**Variety of task**   We analyzed the internal reasoning patterns of language models using synthetic arithmetic reasoning tasks. The use of synthetic data allows for more detailed control compared to experiments on natural language tasks. However, vocabulary and expression diversity, for example, are limited compared to natural language tasks. Therefore, conducting similar analyses on reasoning tasks will verify whether the results of this study apply to other broader, realistic contexts as well.

**Probing methods**   Interpreting internal mechanisms of LMs using probing have been actively conducted in our fields (Conneau et al., 2018; Tenney et al., 2019; Campbell et al., 2023; Li et al., 2023); however, there are criticisms regarding the validity of some probing approaches (Liu et al., 2023; Burns et al., 2023). One way to overcome such concerns will be to analyze the generality of obtained results through more diverse methodologies (Gurnee et al., 2023; Bricken et al., 2023).

## Ethics statement

This paper will not raise particular ethical concerns, considering that (i) no human experiments were conducted, and (ii) our tasks do not involve ethically sensitive topics.

## Acknowledgements

## References

Guillaume Alain and Yoshua Bengio. 2017. Discovering latent knowledge in language models without supervision. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Workshop Track Proceedings*. OpenReview.net.

Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nicholas L. Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E. Burke, Tristan Hume, Shan Carter, Tom Henighan, and Chris Olah. 2023. Towards monosemanticity: Decomposing language models with dictionary learning. In *Anthropic*.

Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. 2023. Discovering latent knowledge in language models without supervision. In *The Eleventh International Conference on Learning Representations*.

James Campbell, Phillip Guo, and Richard Ren. 2023. Localizing lying in llama: Understanding instructed dishonesty on true-false questions through prompting, probing, and patching. In *Socially Responsible Language Modelling Research*.

Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single $&!#* vector: Probing sentence embeddings for linguistic properties. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.

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, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,

Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. The llama 3 herd of models. *CoRR*, abs/2407.21783.

Deep Ganguli, Amanda Askell, Nicholas Schiefer, Thomas I Liao, Kamilė Lukošiūtė, Anna Chen, Anna Goldie, Azalia Mirhoseini, Catherine Olsson, Danny Hernandez, et al. 2023. The capacity for moral self-correction in large language models. *arXiv preprint arXiv:2302.07459*.

Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. Finding neurons in a haystack: Case studies with sparse probing. *Transactions on Machine Learning Research*.

Benjamin Heinzerling and Kentaro Inui. 2024. Monotonic representation of numeric attributes in language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 175–195, Bangkok, Thailand. Association for Computational Linguistics.

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Keito Kudo, Yoichi Aoki, Tatsuki Kuribayashi, Ana Brassard, Masashi Yoshikawa, Keisuke Sakaguchi, and Kentaro Inui. 2023. Do Deep Neural Networks Capture Compositionality in Arithmetic Reasoning? In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 1351–1362.

Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. Emergent World Representations: Exploring a Sequence Model Trained on a Synthetic Task. In *The Eleventh International Conference on Learning Representations (ICLR)*.

Tom Lieberum, Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant Varma, János Kramár, Anca D. Dragan, Rohin Shah, and Neel Nanda. 2024. Gemma scope: Open sparse autoencoders everywhere all at once on gemma 2. *CoRR*, abs/2408.05147.

Kevin Liu, Stephen Casper, Dylan Hadfield-Menell, and Jacob Andreas. 2023. Cognitive dissonance: Why do language model outputs disagree with internal representations of truthfulness? In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4791–4797, Singapore. Association for Computational Linguistics.

Ilya Loshchilov and Frank Hutter. 2017. SGDR: stochastic gradient descent with warm restarts. In *5th International Conference on Learning Representations, ICLR, Toulon, France, April 24-26, 2017, Conference Track Proceedings*.

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. In *Advances in Neural Information Processing Systems*, volume 35, pages 17359–17372. Curran Associates, Inc.

Mistral AI Team. 2024. Mistral nemo.

Timothy Niven and Hung-Yu Kao. 2019. Probing neural network comprehension of natural language arguments. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4658–4664, Florence, Italy. Association for Computational Linguistics.

nostalgebraist. 2020. interpreting gpt: the logit lens. *LessWrong*.

Qwen Team. 2024. Qwen2.5: A party of foundation models.

Herbert E. Robbins. 1951. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407.

Alessandro Stolfo, Yonatan Belinkov, and Mrinmaya Sachan. 2023. A mechanistic interpretation of arithmetic reasoning in language models using causal mediation analysis. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7035–7052, Singapore. Association for Computational Linguistics.

Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Sam Bowman, Dipanjan Das, and Ellie Pavlick. 2019. What do you learn from context? probing for sentence structure in contextualized word representations. In *International Conference on Learning Representations*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. In *Advances in Neural Information Processing Systems 31 (NIPS)*, pages 5998–6008.

Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. Investigating gender bias in language models using causal mediation analysis. In *Advances in Neural Information Processing Systems*, volume 33, pages 12388–12401. Curran Associates, Inc.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le

Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Dingjun Wu, Jing Zhang, and Xinmei Huang. 2023. Chain of Thought Prompting Elicits Knowledge Augmentation. In *Findings of the Association for Computational Linguistics (ACL)*, pages 6519–6534.

An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024a. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*.

Nakyeong Yang, Taegwan Kang, Stanley Jungkyu Choi, Honglak Lee, and Kyomin Jung. 2024b. Mitigating biases for instruction-following language models via bias neurons elimination. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9061–9073, Bangkok, Thailand. Association for Computational Linguistics.

Sohee Yang, Elena Gribovskaya, Nora Kassner, Mor Geva, and Sebastian Riedel. 2024c. Do large language models latently perform multi-hop reasoning? In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10210–10229, Bangkok, Thailand. Association for Computational Linguistics.

Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. 2024. Yi: Open foundation models by 01.ai. *CoRR*, abs/2403.04652.

Fred Zhang and Neel Nanda. 2024. Towards best practices of activation patching in language models: Metrics and methods. In *The Twelfth International Conference on Learning Representations*.

Susan Zhang, Stephen Roller, and Naman Goyal et al. 2022. OPT: Open Pre-trained Transformer Language Models. *CoRR*, arXiv preprint, cs.CL/arXiv:2205.01068.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Ryan Cotterell, Vicente Ordonez, and Kai-Wei Chang. 2019. Gender bias in contextualized word embeddings. 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 629–634, Minneapolis, Minnesota. Association for Computational Linguistics.

| | |
|---|---|
| #train instances | 10,000 |
| Optimizer | SGD (Robbins, 1951) |
| Learning rate | $1.0 \times 10^{-3}$ (constant) |
| Batch size | 10,000 |
| Epochs | 10,000 |

Table 4: Hyperparameters for training the probe

## A  Details of arithmetic tasks

The numbers appearing in each instance are from 0 to 9, and the variable names are sampled from the 26 different alphabets [a-z]. The operators are either addition ($+$) or subtraction ($-$). To standardize the INPUT and OUTPUT token numbers for all evaluation instances at each level, we controlled the equations so that single-digit numbers [0-9] (even after arithmetic operations) were assigned to each variable.

## B  Training settings for linear probes

We trained linear probes in § 2.2. For each combination of input token $t$, layer $l$, and variable $v_i$ in the problem, we train an independent probe $f_{t,l,v_i}(\cdot)$ that predicts the number assigned to $v_i$, given a hidden state $\boldsymbol{h}_{t,l}$:

$$
\begin{aligned}
\hat{y}_{t,l,v_i} &= f_{t,l,v_i}(\boldsymbol{h}_{t,l}) \\
&= \arg\max_C W_{t,l,v_i}\boldsymbol{h}_{t,l} \quad ,
\end{aligned}
\tag{3}
$$

where, $W_{t,l,v_i} \in \mathbb{R}^{|C| \times d}$ is the parameter of the probe. Notably, we formulate the probe as a ten-class classification problem (i.e., $C$=[0-9]) rather than a regression model to align our setting with typical probing experiments. That is, the probe for each combination of $(t, l, v_i)$ is trained to minimize the following cross-entropy loss:

$$
L_{t,l,v_i} = -\sum_{j=1}^{|X|} Y_{v_i}^{[j]} \log(\text{soft}\max_C W_{t,l,v_i}\boldsymbol{h}_{t,l}^{[j]}) \quad ,
\tag{4}
$$

where $|X|$ is the number of test instances, $Y_{v_i}^{[j]} \in \mathbb{R}^C$ is the one-hot encoding of the gold label of a variable $v_i$ in $j$-th test instance, and $\boldsymbol{h}_{t,l}^{[j]}$ is the corresponding hidden state of an LLM in $j$-th instance. The hyperparameters are listed in Table 4.

## C  Additional experiments for § 3

### C.1  Prompt differences

We examined whether altering the method of few-shot prompting influences the model's internal reasoning pattern. Specifically, we investigated the

| Architecture | |
|---|---|
| Base Architecture | OPT (Zhang et al., 2022) |
| Transformer layers | 8 |
| Hidden size | 1,024 |
| Feed-forward size | 4,096 |
| Attention heads | 16 |

| **Hyperparameters for Expert model** | |
|---|---|
| #train instances | 200,000 |
| Optimizer | Adam (Kingma and Ba, 2015) ($\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 1 \times 10^{-8}$) |
| Learning rate scheduler | cosine (Loshchilov and Hutter, 2017) |
| Learning rate warmup | 500 |
| Learning rate (maximum) | $3.0 \times 10^{-5}$ |
| Gradient clipping | 1.0 |
| Dropout | 0.1 |
| Batch size | 16,384 tokens |
| Epochs | 500 |

| **Hyperparameters for General model** | |
|---|---|
| #train instances | 1,000,000 examples |
| Optimizer | Adam (Kingma and Ba, 2015) ($\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 1 \times 10^{-8}$) |
| Learning rate scheduler | cosine (Loshchilov and Hutter, 2017) |
| Learning rate warmup | 500 |
| Learning rate (maximum) | $1.0 \times 10^{-4}$ |
| Gradient clipping | 1.0 |
| Dropout | 0.1 |
| Batch size | 16,384 tokens |
| Epochs | 50 |

Table 5: Hyperparameters and model architecture for General Model and Expert model

differences in internal reasoning patterns between cases where the model is provided with three equations of the same level as the evaluation task (same level prompting) and cases where the model is provided with 50[2] randomly generated equations consisting of three equations (general prompting), as described in Section 2.1. Table 7 shows a comparison of RSP between the two prompting methods. The experimental results showed no difference in RSP between these two prompting methods. Therefore, it is considered that the impact of the format of few-shot prompting on the internal reasoning pattern is limited.

### C.2  Models trained from scratch

As a sanity check to confirm that our probing approach indeed renders different results depending on the targeted models, we compare the probing results between two models trained on extremely different data.

---

[2]Due to the insufficient performance of the model on the task, the number of demonstrations in general prompts is larger than that in same level prompts.

|  |  | RSP (↓) | | Acc (↑) | |
|  | | #steps | eq | tok | Max | Task |
|---|---|---|---|---|---|---|
| Expert model | $v_0$ | 2 | $-2$ | $-8$ | 100.0 | |
| | $v_1$ | 1 | - | - | 83.2 | 100.0 |
| | $v_2$ | 1 | - | - | 21.5 | |
| General model | $v_0$ | 2 | $-2$ | $-8$ | 100.0 | |
| | $v_1$ | 1 | $-2$ | $-8$ | 98.8 | 99.9 |
| | $v_2$ | 1 | $-1$ | $-2$ | 98.2 | |

Table 6: The results of the Expert model and General model. The "eq" columns indicate the model's $\text{RSP}_{\text{eq}}$ score. The "tok" column shows the $\text{RSP}_{\text{tok}}$ score The Max scores refers to the highest accuracy achieved by the probe. Task represents the accuracy of the model in the evaluation task. "-" indicates that the accuracy of the probes never reached the threshold ($\tau = 0.9$).

### C.2.1 Settings

Specifically, we analyze the two transformer-based LMs [3]:

- EXPERT model: This model is trained only on level 4. That is, the model can know, for example, C ($v_2$) is unnecessary to calculate the final answer A=? ($v_0$) and memorize the solution, i,e., add the first three numbers.

- GENERAL model: This model is trained on problems consisting of random combinations of three or fewer equations.

The texts are tokenized at the character level. We use the absolute position embeddings from the Transformer (Vaswani et al., 2017) as embeddings for each digit to make it easy for the model to learn the magnitude relationships between digits. We use randomly initialized embeddings for variables, operators, and other non-numeric embeddings, and only embeddings are frozen during training. Other hyperparameters and model architectures are listed in Table 5.

### C.2.2 Experimental results

We employed the level 4 task introduced in Table 1. Figure 3 and Table 6 show the results of the probing for the GENERAL model and EXPERT model. The results show strikingly different pictures between the two models. For example, in the EXPERT model, the probe accuracy for B ($v_1$) never reaches nearly 100% (Table 6) probably due to its memorization of shallow solution (i.e., add the first three numbers). In contrast, in the GENERAL model,

---

[3] We used Hugging Face Transformers (Wolf et al., 2020) as the implementation.

| Prompting fromat | Variable name | #steps | RSP (↓) eq* | eq | Acc (↑) < CoT | > CoT | Task |
|---|---|---|---|---|---|---|---|
| same level | $v_0$ | 2 | $-3$ | 4 | 46.0 | 100 | |
| | $v_1$ | 1 | $-3$ | $-3$ | 99.2 | 100 | 100 |
| | $v_2$ | 1 | $-2$ | $-2$ | 99.1 | 34.5 | |
| general | $v_0$ | 2 | $-3$ | 4 | 45.8 | 99.9 | |
| | $v_1$ | 1 | $-3$ | $-3$ | 99.9 | 100 | 99.9 |
| | $v_2$ | 1 | $-2$ | $-2$ | 99.8 | 45.6 | |

Table 7: Evaluation results of same level promting and general prompting. The eq* and eq columns indicate the lower bound and model's $\text{RSP}_{\text{eq}}$ score. The $<$ CoT and $>$ CoT scores correspond to the accuracies introduced in § 2.3. Task represents the accuracy of the model in the evaluation task.

the probe accuracy for A ($v_0$) increases, followed by an increase in the probe accuracy for B ($v_1$) (Figure 3).

Relatedly, the distractor equation ($v_2 = 4 + 5$) is only solved by the GENERAL model, i.e., the EXPERT model skipped the calculation of the distractor. Specifically, the probe accuracy for predicting the value of $v_2$ reaches a maximum of 98.2% in the GENERAL model, while it is only 21.5% in the EXPERT model. This suggests that the GENERAL model performs calculations for all equations, regardless of the necessity for the final answer. Such different and intuitive results between the two models support the fact that our probing approach can indeed render their different internals for different models.

Notably, one can also see some consistent patterns between the two models; for example, in both EXPERT and GENERAL models, the probe accuracy for predicting B ($v_1$) increases immediately after they take necessary information to determine the value of B ($v_1$).

## D   Discussion of results in § 3

**Distractors:** In task level 4, the value of C ($v_2$) is unnecessary for determining the final answer (distractor), and this fact can be guessed based on the in-context examples. From Table 2, the maximum accuracy of the probe for C $v_2$ in level 4 was almost 100%. Furthermore, the $\text{RSP}_{\text{tok}}(v_2)$ matched the lower bound (eq*). This suggests that the models tend to calculate equations even when they are unnecessary for the final answer, implying their somewhat redundant internal reasoning.

**Error analysis:** Figure 4 presents the top-1 predictions of the probe for B on instances where Llama3.2-3B generated incorrect answers for Task
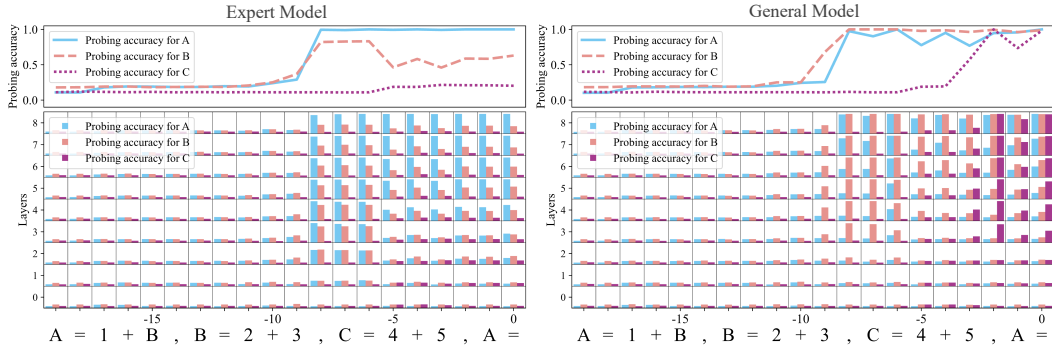
Figure 3: Probing results for models trained only on arithmetic reasoning tasks. The vertical axis represents the index of the transformer layer. The horizontal axis represents the tokens input to the model over time. The vertical axis of the line graph shows the maximum accuracy of the probes at each position (time) $t$. It can be observed that the probe accuracy for the value of $v_2$, an unnecessary equation for the final answer, is low in the Expert model.

3）．The equations at the bottom represent the answers generated by the model, with the incorrect responses highlighted in red. Figure 4 illustrates that the model successfully derives the correct answer during inference (green text in Figure 4), but it transitions to an incorrect answer immediately before output (red text in Figure 4). A detailed analysis of what causes this shift to an incorrect answer just before output is left for future work.

## E    Detailes of causal interventions § 4

In § 4, we employed activation patching to validate the probing results and explore the underlying mechanisms to generate the final answer. Here, we provide details of the experimental setup and the results of probing when activation patching is applied.

**Experimental setup**    As described in Section 4.1, we applied activation patching to the hidden states with high probe accuracy ($H^+$) and low probe accuracy ($H^-$). Our analysis focused on the Qwen2.5-7B at task level 3. Based on the probe accuracy for the variable B ($v_1$) shown in Figure 2, we defined $H^+$ and $H^-$ as follows: $H^+ = \{h_{t,l} \mid l \geq 12, 13 \leq t \leq 14\}$, $H^- = \{h_{t,l} \mid l \geq 12, 15 \leq t \leq 17\}$.

**Experimental results**    Table 8 aggregates the results of our causal intervention analysis.

**Visualization of activation patching**    Figure 5 shows the visualization of the probing results when activation patching is applied. The gray boxes highlight the areas where the hidden states of other problems have been patched. The accuracy of the probes in the patched areas exhibits a significant

| Setting | Patch states | OUTPUT | Unchanged rate | Success rate |
|---------|--------------|--------|----------------|--------------|
| 1 | $H^+$ | self | 0.00 | 100 |
| 2 | $H^-$ | self | 100 | 0.00 |
| 3 | $H^+$ | original | 100 | 0.00 |

Table 8: OUTPUT refers to the configuration of the language model's output. The "self" setting allows the model to generate text by the model itself. The "original" setting refers to the force-decoding of gold labels when patching is not applied. The unchanged rate represents the ratio of final answers that remain constant despite activation patching. The success rate indicates the ratio of final answers that changed to a new output upon applying a patch state.

decline, confirming that the model's internal state changed.

## F    All probing results

Figures 6 to 18 show the probing results for all models and tasks mentioned in this paper.

## G    Computational resources

We used NVIDIA A100 GPUs (40GB and 80GB memory) to conduct this study.

## H    Usage of AI assistants

For writing this paper and the source code for the experiments, we use AI assistants (e.g., ChatGPT, GitHub Copilot). However, the use is limited to purposes such as code completion, translation, text editing, and table creation, and the content is based on the authors' ideas.
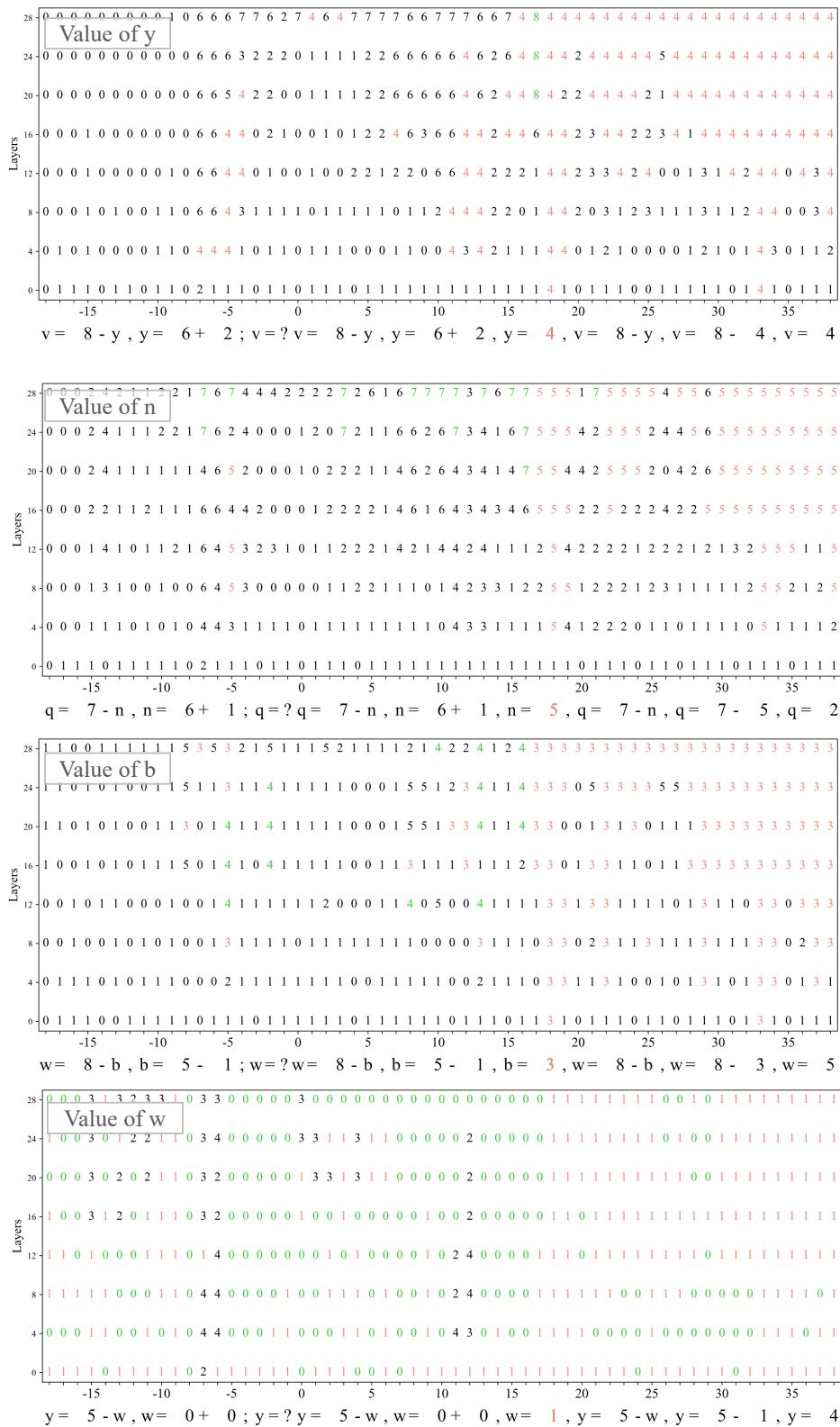
Figure 4: Error analysis of cases where Llama3.2-3B generated incorrect answers. The vertical axis represents the index of the transformer layer. The horizontal axis represents the tokens input to the model over time. The numbers in the figure indicate the labels predicted by each probe (top-1) at each layer and time step. The numbers highlighted in green represent the gold labels for the predictions, while those highlighted in red denote the values incorrectly generated by the language model.

Figure 5: Probing results when activation patching is applied. The vertical axis represents the index of the transformer layer. The horizontal axis represents the tokens input to the model over time. The gray boxes highlight the areas where the hidden states patched.



Figure 6: Probing results when Qwen2.5-7B solves Task 1.

Figure 7: Probing results when Qwen2.5-7B solves Task 2.
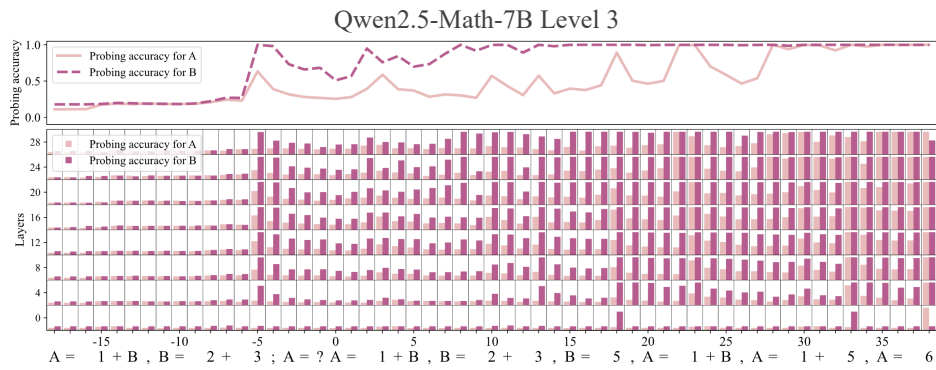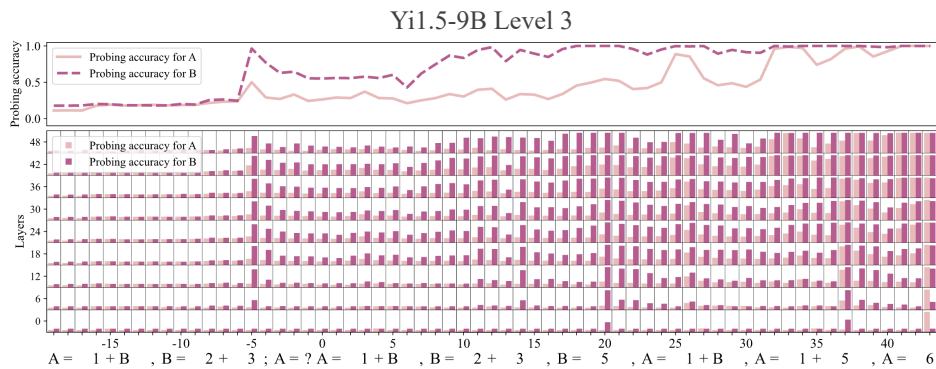


Figure 8: Probing results when Qwen2.5-7B solves Task 3.



Figure 9: Probing results when Qwen2.5-7B solves Task 4.



Figure 10: Probing results when Qwen2.5-7B solves Task 5.
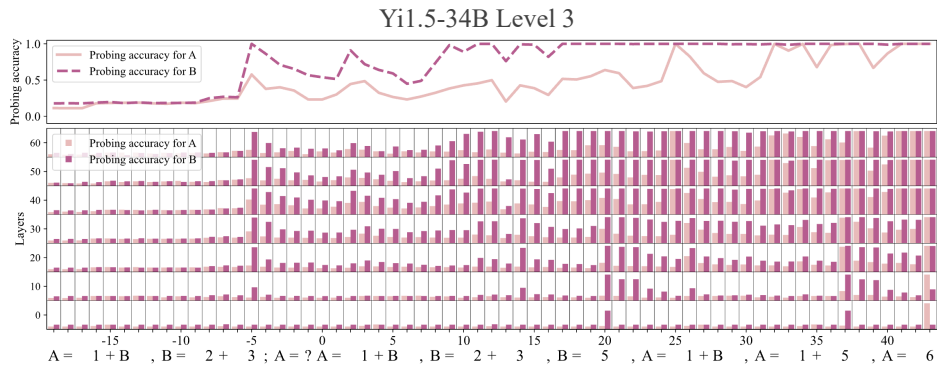
Figure 11: Probing results when Qwen2.5-14B solves Task 3.



Figure 12: Probing results when Qwen2.5-32B solves Task 3.



Figure 13: Probing results when Qwen2.5-Math-7B solves Task 3.



Figure 14: Probing results when Yi1.5-9B solves Task 3.

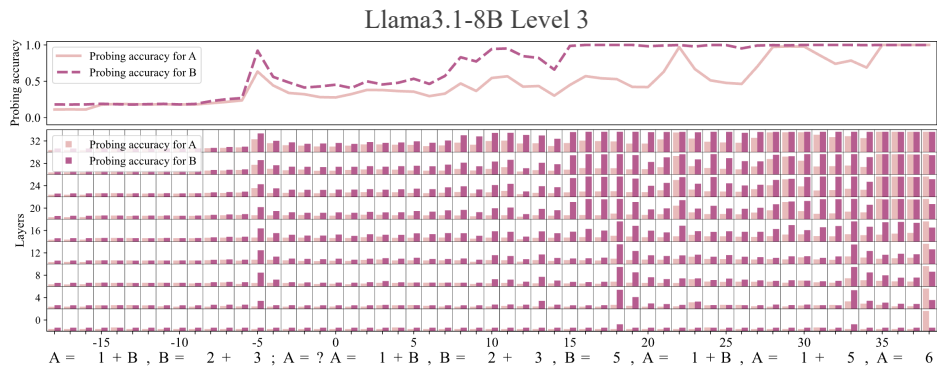Figure 15: Probing results when Yi1.5-34B solves Task 3.



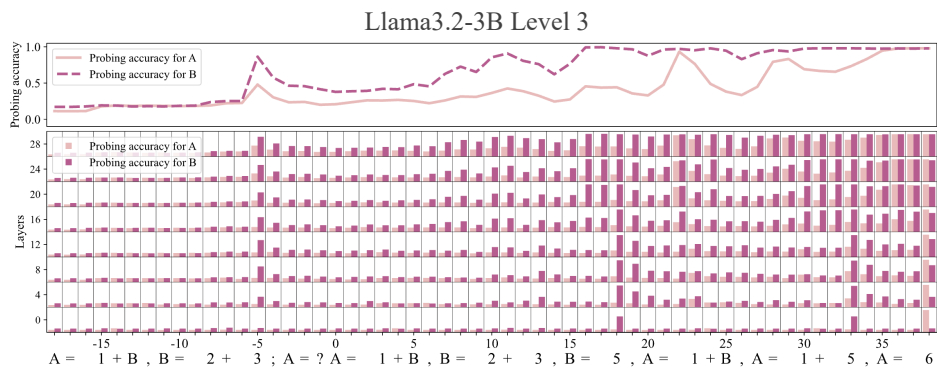Figure 16: Probing results when Llama3.1-8B solves Task 3.
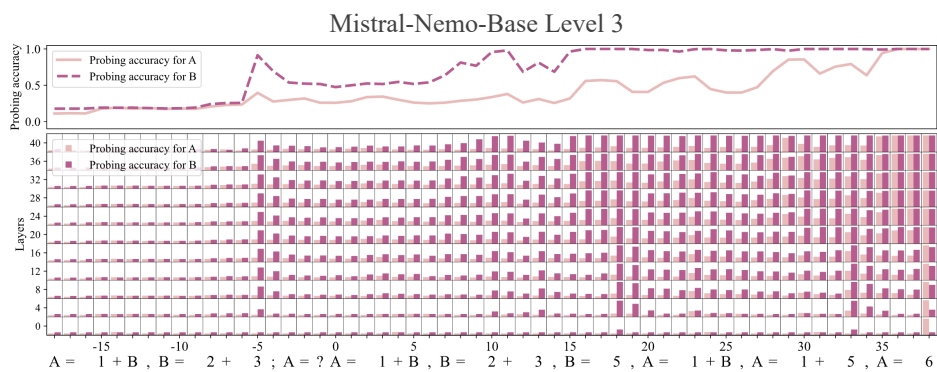


Figure 17: Probing results when Llama3.2-3B solves Task 3.



Figure 18: Probing results when Mistral-Nemo-Base solves Task 3.