

Are Reflective Words in Large Reasoning Models a Sign of Genuine Capability or Memorized Patterns?

Anonymous ACL submission

Abstract

Recent Large Reasoning Models exhibit strong reasoning abilities in tasks like mathematics and logical inference, notably through human-like self-verification and reflection in their chain of thought. However, it remains unclear whether these reflective statements stem from genuine internal mechanisms or are merely memorized patterns. From a model interpretability perspective, this work investigates LRMs’ representation space to determine whether specific features causally govern reflective capabilities. Using a difference-in-means approach, we extract Self-Reflection Features by contrasting model activations during self-reflection versus affirmative answering. Further causal analysis reveals that these features strongly influence knowledge parameters associated with reflection words, suggesting that such outputs are genuine manifestations of internal mechanisms rather than memorization. Finally, causal interventions demonstrate that modulating these features flexibly adjusts the model’s self-reflective intensity.

1 Introduction

Recently, the emergence of Large Reasoning Models (LRMs) (OpenAI et al., 2024b; DeepSeek-AI et al., 2025; Team, 2025) optimized through reinforcement learning, has opened up new possibilities and room for advancement in the reasoning capabilities of language models. These advancements are particularly evident in tasks such as mathematics (Cobbe et al., 2021; Hendrycks et al., 2021), logical reasoning (Luo et al., 2024), and understanding scientific questions (Welbl et al., 2017). These Reasoning models excel at deconstructing complex problems into simpler, sequential sub-problems within their extensive chains of thought. Most impressively, they often adopt a human-like reasoning tone (Guo et al., 2025; Yang et al., 2025), seemingly engaging in self-verification and reflection (Gandhi et al., 2025), and evaluating their own

proposed solutions before summarizing and recommending the most suitable option to the user.

Therefore, this raises a crucial question: do these reflective statements and verification words executed within the chain of thought represent **a genuine activation of the models’ internal reflective capabilities, or are they simply reproductions of patterns memorized from their training data?**

In this work, from a model interpretability perspective, we delve into the representation space of LRMs. We aim to uncover whether specific existing features genuinely govern the deployment of these reflective capabilities, and to establish if a causal relationship exists between the activation of such features and the reflection words manifested in a reasoning model’s chain of thought.

Specifically, in §3, we employ the *difference-in-means* technique (Marks and Tegmark, 2023; Rimsky et al., 2024) to extract **Self-Reflection Features** from four Large Reasoning Models (Guo et al., 2025; Team, 2025) across both mathematical and code datasets by contrasting the internal representations of the models when they engage in self-reflection versus when they provide affirmative answers. These features were subsequently visualized using Principal Component Analysis.

In §4, we conduct a causal analysis of Self-Reflection Features in LRMs from both internal and external perspectives. We identify knowledge parameters within the models that are highly correlated with reflection words and demonstrate that the presence of Self-Reflection Features amplifies the activation of these parameters in §4.1. This suggests that the reflection words in LRM chain-of-thought are genuine manifestations of these activated features, not just memorized patterns. Crucially, through causal intervention experiments detailed in §4.2, we further show that manipulating these extracted Self-Reflection Features allows for flexible modulation of the model’s self-reflection intensity when answering questions. To conclude,

we uncovered and verified that the reflection words in Large Reasoning Models genuinely reflect the activation of their internal reflective capabilities.

2 Background and Related Work

2.1 Self-Reflection in Large Reasoning Models

The development of Large Reasoning Models (LRMs) (OpenAI et al., 2024b; Guo et al., 2025; Team, 2025) has opened up new prospects for enhancing the reasoning paradigms of language models. Most notably, they demonstrate impressive human-like self-reflection (Guo et al., 2025; Liu et al., 2025) and verification capabilities when engaged in the long chain of thoughts (Wei et al., 2023; Li et al., 2025).

Regarding the **human-like expressions** in the chain of thoughts (CoT) exhibited by LRMs, several studies have conducted preliminary investigations from the perspective of Reinforcement Learning training dynamics (Gandhi et al., 2025; Yang et al., 2025; Yu et al., 2025b). And in terms of the LRM’s ability to **assess itself’s uncertainty** or engage in self-reflection, existing research has explored both explicit and implicit ways to estimate the uncertainty. For explicit uncertainty, prior work proposed prompting strategies that guide LRMs to verbalize their confidence levels (Zeng et al., 2025). To study implicit uncertainty, researchers have trained probing classifiers on the model’s internal representations to estimate its confidence (Zhang et al., 2025; Anthropic, 2025). However, there is still a lack of sufficient interpretability research exploring whether these explicit reflection patterns observed in COT genuinely correlate with the models’ actual internal reflective capabilities.

2.2 Linear semantic features

Recent investigations in model interpretability have revealed that, for numerous cognitive behaviors observed in Large Language Models—including refusal to answer (Arditi et al., 2024), jailbreaking (Yu et al., 2025a), reasoning, and knowledge-recall (Hong et al., 2025)—the models encode corresponding linear semantic features within their activation space (Park et al., 2024). These linear semantic features have been discovered and extracted by contrasting inputs that differ primarily in the target semantic dimension (Marks and Tegmark, 2023). Once these features are pinpointed, they offer a mechanism for controlling model behavior through manipulation, which allows for targeted in-

terventions in the generative process (Rimsky et al., 2024; Stickland et al., 2024). Our work extends this line of study by identifying linear features that determine models’ engagement in self-reflection.

3 Self-Reflection Features Extraction

3.1 Methodology for Identifying Self-Reflection Features

For current Reasoning Models, given a question Q , we can decompose its output into multiple *Reasoning Segments*: $\{s_1, s_2, s_3, \dots, s_n\}$. Each segment (except for s_1) represents the model’s reflection on the previous segment’s proposed approach and a new attempt at solving the target problem. The final segment, s_n , signifies the termination of reflection, and the model directly provides its final answer.

At each *Reasoning Segment*’s final token position during inference, the model can either select the current answer as its final output and terminate, or generate another segment to reflect, verify, and explore alternative solutions. Therefore, based on whether the model initiates a new reflection after a segment or directly provides the final answer, we can categorize the *Reasoning Segments* into two groups. The first group, where the model proposes a new reflection after the segment, we call $\mathcal{S}_{\text{Check-point}}$. The second group, where the model directly gives the final answer after the segment, we call $\mathcal{S}_{\text{Termination}}$.

Next, for both groups, we extract the hidden states from the last-token position of each segment s (excluding s_n) at the model’s l -th layer¹, denoted as $h^{(l)}(s)$. Since this last-token position corresponds to where the model is about to generate the first token of the subsequent segment, we hypothesize that the hidden states at this crucial juncture store important information guiding the model’s decision to either continue with reflection or proceed to termination in the next segment. Then, using the *difference-in-means* technique (Marks and Tegmark, 2023; Rimsky et al., 2024), we calculate the difference between the mean last-token hidden states for these two categories of *Reasoning Segments*:

$$\mathbf{f}^{(l)} = \frac{\sum_{s \in \mathcal{S}_{\text{Check-point}}} \mathbf{h}^{(l)}(s)}{|\mathcal{S}_{\text{Check-point}}|} - \frac{\sum_{s \in \mathcal{S}_{\text{Termination}}} \mathbf{h}^{(l)}(s)}{|\mathcal{S}_{\text{Termination}}|} \quad (1)$$

¹Assuming the model has L layers, we conduct experiments on each individual layer of it.

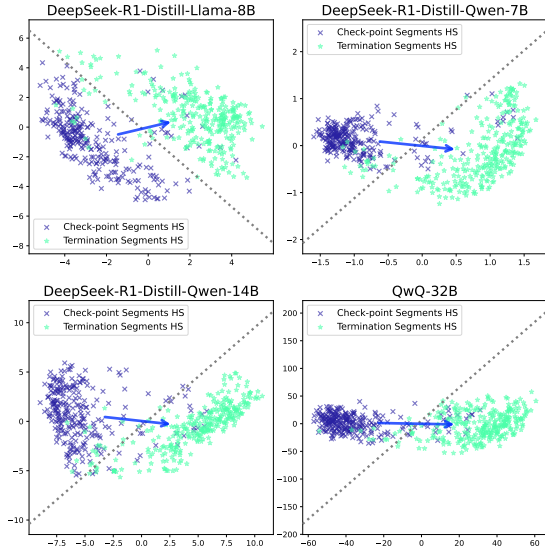


Figure 1: Visualization of the hidden states of four reasoning models on GSM8k dataset using 2-D PCA. The hidden states of datapoints in $\mathcal{S}_{Check-point}$ and $\mathcal{S}_{Termination}$ are positioned around the boundary (grey dashed line) fitted via logistic regression. The blue arrow approximately indicates the direction of the Self-Reflection Features. Results on other datasets are shown in §C.

The direction of the vector $\mathbf{f}^{(l)}$ represents the direction of the Self-Reflection Features that we extracted. The construction details of $\mathcal{S}_{Check-point}$ and $\mathcal{S}_{Termination}$ are provided in the next section.

3.2 Experimental Setups

Models We utilize two categories of reasoning models trained under different settings to investigate Self-Reflection Features. The first category includes DeepSeek-R1-Distill-Llama-8B, DeepSeek-R1-Distill-Qwen-7B, and DeepSeek-R1-Distill-Qwen-14B. These models are obtained by performing supervised fine-tuning on the base Llama-3.1 (Meta, 2024) or Qwen2.5 (Qwen et al., 2025) models using high-quality reasoning data generated by the DeepSeek-R1 model (Guo et al., 2025). The second category comprises the QwQ-32B model (Team, 2025), which is trained using reinforcement learning. The inference details are provided in §E.

Datasets We focus on analyzing LRMs on mathematical and coding tasks to facilitate the extraction of Self-Reflection Features and analyze their influence. For the mathematical tasks, we use GSM8k (Cobbe et al., 2021) and MATH-500 datasets (Lightman et al., 2023). For the coding tasks, we select MBPP dataset (Austin et al., 2021).

3.3 Visualization for Self-Reflection Features

Following the methodology outlined in §3.1, we first perform inference on the datasets using the target reasoning models to collect their responses. We then employ GPT-4o (OpenAI et al., 2024a) to automatically segment each response into multiple reasoning segments², where each segment independently represents an attempt by the model to solve the problem. Subsequently, based on the segmentation results, we categorize the segments into two groups, $\mathcal{S}_{Check-point}$ and $\mathcal{S}_{Termination}$. Detailed statistics for these two groups are presented in §B.

We then extract the hidden states from the corresponding positions and compute the Self-Reflection Features by applying Eq. (1). To more clearly visualize the direction of the Self-Reflection Features, we apply Principal Component Analysis (PCA) to the hidden states of data points in the $\mathcal{S}_{Check-point}$ and $\mathcal{S}_{Termination}$ sets. The results on GSM8k dataset are shown in Figure 1. From this, we observe that two groups of data points are clearly divisible into two clusters by a logistic regression line, explicitly revealing the presence of Self-Reflection Features.

4 Internal and External Causal Analysis of Self-Reflection Features in LRMs

In this section, we will investigate from both internal (model parameter activation) and external (runtime behavior) perspectives, to verify the genuine causal relationships connecting Self-Reflection Features with: (a) the presence of reflection words within chain-of-thought processes in §4.1, and (b) the intensity of the model’s self-reflection during actual inference in §4.2.

4.1 Parameter Storing Human-like Reflection words

By applying the Logit Lens method (nostalgebraist, 2020)³, we identified value vectors within the MLP module’s value matrix of the large reasoning models that highly contain these reflection tokens. Specific examples are presented in Table 1. We can observe that the vector projections at corresponding positions in both models each contain a certain number of reflection tokens. Moreover, when the hidden states are in $\mathcal{S}_{Check-point}$ — that is, when they exhibit stronger self-reflection features — we ob-

²The exact prompts used, along with human verification results, are provided in §A of the appendix.

³More descriptions of this method and relevant background knowledge are provided in §D of the appendix.

Model	Example Top-scoring tokens Vector	Activation values in $S_{Termination}$	Activation values in $S_{Check-point}$	Activation values in $S_{More-Reflection}$
DeepSeek-R1-Distill-Llama-8B	\mathbf{v}_{10644}^{31} Is, Let, OK, So, Next, If, Now, What, First, See, However, Like, Check, Right, Wait, Again	0.15	2.13 $\uparrow 2.0$	2.71 $\uparrow 2.6$
DeepSeek-R1-Distill-Qwen-7B	\mathbf{v}_{11862}^{23} hi, well, its, hey, nah, its, Im, ye, alternative, _ok,, oh, Hello, notifies, thanks, Ye, waits, WAIT	0.22	1.94 $\uparrow 1.7$	2.40 $\uparrow 2.2$

Table 1: Example value vectors identified in two LRMs via the Logit Lens method, showcasing top-scoring tokens related to Self-Reflection Features and their activation values in different reflection stages.

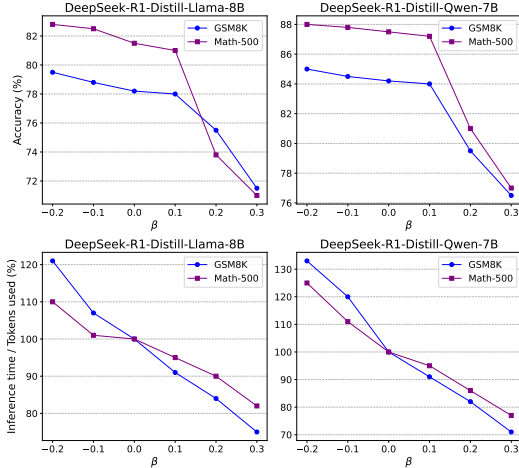


Figure 2: Accuracy and inference time on GSM8K and MATH-500 after intervening in the hidden states of DeepSeek-R1-Distill-Llama-8B and DeepSeek-R1-Distill-Qwen-7B using different β values in Eq. (2).

serve a significant increase in the activation⁴ of the target value vector. When we follow Eq. (2) to further enhance the strength of self-reflection features in the model representations (i.e., transitioning to $S_{More-Reflection}$ as shown in Table 2), we similarly observe a further increase in their activation.

This provides further support for the idea that the human-like reflection words appearing in the chain-of-thought processes of LRMs are not merely a result of memorizing training data. Instead, they are **a reflection of the genuine activation of Self-Reflection Features**.

4.2 Modulating Self-Reflection Intensity via Linear Feature Intervention

Building upon the Linear Reflection Features extracted from LRMs in §3, in this part we aim to modulate the intensity of model’s self-reflection features within the model’s representational space during inference for specific tasks, thereby address-

⁴Activation refers to the coefficient corresponding to each value vector in Eq. (4).

ing the potential issues of insufficient reflection (Aggarwal and Welleck, 2025) or "overthinking" (Cuadron et al., 2025; Zhang et al., 2025) that current LRMs may exhibit in practical scenarios.

Specifically, following Eq. (2) below, we attempt to control the extent of Self-Reflection ability in LRMs by intervening with Self-Reflection Features through adjusting the hyperparameter β in the model’s hidden states:

$$\mathbf{h}'^{(l)}(s) \leftarrow \mathbf{h}^{(l)}(s) - \beta * \mathbf{f}^{(l)}, \quad (2)$$

where $\mathbf{h}^{(l)}(s)$ and $\mathbf{h}'^{(l)}(s)$ denote the hidden states in the l -th layer of the model before and after intervention, respectively, when performing inference on segment s . And $\mathbf{f}^{(l)}$ represents the feature direction we extracted using Eq. 1. Figure 2 shows the intervention effects on GSM8K and MATH-500 for both DeepSeek-R1-Distill-Llama-8B and DeepSeek-R1-Distill-Qwen-7B. Starting from $\beta = 0$, increasing the value of β leads to improved acceleration for the LRMs. However, accuracy does not immediately degrade. Once β reaches around 0.2, further acceleration comes at the cost of a noticeable drop in accuracy. Conversely, decreasing β —thereby increasing the influence of Self-Reflection Features—slightly improves accuracy at the expense of slower inference, which aligns with our hypothesis.

5 Discussion and Conclusion

This work aimed to determine if reflective language in LRMs reflects genuine internal processes or learned patterns. By extracting Self-Reflection Features from the representation space, we found a causal link to reflective words in their chain of thought. Crucially, manipulating these features allowed us to modulate LRM self-reflection intensity. These findings confirm LRMs’ reflective statements stem from discernible, governable internal mechanisms, signifying true reflective activation.

305 Limitations

306 In this study, while we have identified the pres-
307 ence of Self-Reflection features within reasoning
308 models, a comprehensive investigation into their
309 **origins** was not conducted. Specifically, it remains
310 to be clarified whether these features emerge pri-
311 marily from the pre-training phase or are intro-
312 duced during subsequent reinforcement learning
313 post-training. Furthermore, the characteristics of
314 the training data that facilitate the encoding of these
315 Self-Reflection features into the model’s represen-
316 tational space are yet to be identified. A deeper
317 understanding of these aspects would provide a
318 critical foundation for the future development of
319 more robust and effective reasoning models. We
320 plan to explore these questions in our future work.

321 Additionally, owing to resource constraints, we
322 were unable to extend our experimental research to
323 larger-scale reasoning models, such as DeepSeek-
324 R1.

325 References

326 Pranjali Aggarwal and Sean Welleck. 2025. **L1: Con-**
327 **trolling how long a reasoning model thinks with rein-**
328 **forcement learning.** *Preprint*, arXiv:2503.04697.

329 Anthropic. 2025. On the biology of a large language
330 model. [https://transformer-circuits.pub/](https://transformer-circuits.pub/2025/attribution-graphs/biology.html)
331 [2025/attribution-graphs/biology.html](https://transformer-circuits.pub/2025/attribution-graphs/biology.html). Ac-
332 cessed: 2025-03-27.

333 Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka,
334 Nina Panickssery, Wes Gurnee, and Neel Nanda.
335 2024. Refusal in language models is mediated by
336 a single direction. *arXiv preprint arXiv:2406.11717*.

337 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten
338 Bosma, Henryk Michalewski, David Dohan, Ellen
339 Jiang, Carrie Cai, Michael Terry, Quoc Le, and
340 Charles Sutton. 2021. **Program synthesis with large**
341 **language models.** *Preprint*, arXiv:2108.07732.

342 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
343 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
344 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
345 Nakano, Christopher Hesse, and John Schulman.
346 2021. **Training verifiers to solve math word prob-**
347 **lems.** *Preprint*, arXiv:2110.14168.

348 Alejandro Cuadron, Dacheng Li, Wenjie Ma, Xingyao
349 Wang, Yichuan Wang, Siyuan Zhuang, Shu Liu,
350 Luis Gaspar Schroeder, Tian Xia, Huanzhi Mao,
351 Nicholas Thumiger, Aditya Desai, Ion Stoica, Ana
352 Klimovic, Graham Neubig, and Joseph E. Gonzalez.
353 2025. **The danger of overthinking: Examining the**
354 **reasoning-action dilemma in agentic tasks.** *Preprint*,
355 arXiv:2502.08235.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, 356
Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, 357
Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, 358
Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zihong 359
Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, 360
Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, 361
Chenggang Zhao, Chengqi Deng, Chenyu Zhang, 362
Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, 363
Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, 364
Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, 365
Han Bao, Hanwei Xu, Haocheng Wang, Honghui 366
Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, 367
Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang 368
Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. 369
Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai 370
Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai 371
Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong 372
Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan 373
Zhang, Minghua Zhang, Minghui Tang, Meng Li, 374
Miaojun Wang, Mingming Li, Ning Tian, Panpan 375
Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, 376
Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, 377
Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, 378
Shanghai Lu, Shangyan Zhou, Shanhuang Chen, 379
Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng 380
Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing 381
Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, 382
T. Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, 383
Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao 384
Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan 385
Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin 386
Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, 387
Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, 388
Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxi- 389
ang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, 390
Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang 391
Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng 392
Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, 393
Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, 394
Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, 395
Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yu- 396
jia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, 397
Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, 398
Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, 399
Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, 400
Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean 401
Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, 402
Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zi- 403
jia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, 404
Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu 405
Zhang, and Zhen Zhang. 2025. **Deepseek-r1: Incent-**
406 **ivizing reasoning capability in llms via reinforce-**
407 **ment learning.** *Preprint*, arXiv:2501.12948. 408

Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh, 409
Nathan Lile, and Noah D. Goodman. 2025. **Cogni-**
410 **tive behaviors that enable self-improving reasoners,**
411 **or, four habits of highly effective stars.** *Preprint*,
412 arXiv:2503.01307. 413

Mor Geva, Avi Caciularu, Guy Dar, Paul Roit, Shoval 414
Sadde, Micah Shlain, Bar Tamir, and Yoav Goldberg. 415
2022a. **LM-debugger: An interactive tool for inspec-**
416 **tion and intervention in transformer-based language**
417

535	Pachocki, James Aung, James Betker, James Crooks,	Degry, Thomas Dimson, Thomas Raoux, Thomas	599
536	James Lennon, Jamie Kiros, Jan Leike, Jane Park,	Shadwell, Tianhao Zheng, Todd Underwood, Todor	600
537	Jason Kwon, Jason Phang, Jason Teplitz, Jason	Markov, Toki Sherbakov, Tom Rubin, Tom Stasi,	601
538	Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Var-	Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce	602
539	avva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui	Walters, Tyna Eloundou, Valerie Qi, Veit Moeller,	603
540	Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang,	Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne	604
541	Joaquin Quinonero Candela, Joe Beutler, Joe Lan-	Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra,	605
542	ders, Joel Parish, Johannes Heidecke, John Schul-	Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian,	606
543	man, Jonathan Lachman, Jonathan McKay, Jonathan	Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen	607
544	Uesato, Jonathan Ward, Jong Wook Kim, Joost	He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and	608
545	Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross,	Yury Malkov. 2024a. Gpt-4o system card . <i>Preprint</i> ,	609
546	Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao,	arXiv:2410.21276.	610
547	Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai		
548	Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin	OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer,	611
549	Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu,	Adam Richardson, Ahmed El-Kishky, Aiden Low,	612
550	Kenny Nguyen, Keren Gu-Lemberg, Kevin Button,	Alec Helyar, Aleksander Madry, Alex Beutel, Alex	613
551	Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle	Carney, Alex Iftimie, Alex Karpenko, Alex Tachard	614
552	Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lau-	Passos, Alexander Neitz, Alexander Prokofiev,	615
553	ren Workman, Leher Pathak, Leo Chen, Li Jing, Lia	Alexander Wei, Allison Tam, Ally Bennett, Ananya	616
554	Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lil-	Kumar, Andre Saraiva, Andrea Vallone, Andrew Du-	617
555	lian Weng, Lindsay McCallum, Lindsey Held, Long	berstein, Andrew Kondrich, Andrey Mishchenko,	618
556	Ouyang, Louis Fevrier, Lu Zhang, Lukas Kon-	Andy Applebaum, Angela Jiang, Ashvin Nair, Bar-	619
557	draciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz,	ret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin	620
558	Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine	Sokolowsky, Boaz Barak, Bob McGrew, Borys Mi-	621
559	Boyd, Madeleine Thompson, Marat Dukhan, Mark	naiev, Botao Hao, Bowen Baker, Brandon Houghton,	622
560	Chen, Mark Gray, Mark Hudnall, Marvin Zhang,	Brandon McKinzie, Brydon Eastman, Camillo Lu-	623
561	Marwan Aljubeih, Mateusz Litwin, Matthew Zeng,	garesi, Cary Bassin, Cary Hudson, Chak Ming Li,	624
562	Max Johnson, Maya Shetty, Mayank Gupta, Meghan	Charles de Bourcy, Chelsea Voss, Chen Shen, Chong	625
563	Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao	Zhang, Chris Koch, Chris Orsinger, Christopher	626
564	Zhong, Mia Glaese, Mianna Chen, Michael Jan-	Hesse, Claudia Fischer, Clive Chan, Dan Roberts,	627
565	ner, Michael Lampe, Michael Petrov, Michael Wu,	Daniel Kappler, Daniel Levy, Daniel Selsam, David	628
566	Michele Wang, Michelle Fradin, Michelle Pokras,	Dohan, David Farhi, David Mely, David Robinson,	629
567	Miguel Castro, Miguel Oom Temudo de Castro,	Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Free-	630
568	Mikhail Pavlov, Miles Brundage, Miles Wang, Mi-	man, Eddie Zhang, Edmund Wong, Elizabeth Proehl,	631
569	nal Khan, Mira Murati, Mo Bavarian, Molly Lin,	Enoch Cheung, Eric Mitchell, Eric Wallace, Erik	632
570	Murat Yesildal, Nacho Soto, Natalia Gimelshein, Na-	Ritter, Evan Mays, Fan Wang, Felipe Petroski Such,	633
571	talie Cone, Natalie Staudacher, Natalie Summers,	Filippo Raso, Florencia Leoni, Foivos Tsimpourlas,	634
572	Natan LaFontaine, Neil Chowdhury, Nick Ryder,	Francis Song, Fred von Lohmann, Freddie Sulit,	635
573	Nick Stathas, Nick Turley, Nik Tezak, Niko Felix,	Geoff Salmon, Giambattista Parascandolo, Gildas	636
574	Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel	Chabot, Grace Zhao, Greg Brockman, Guillaume	637
575	Bundick, Nora Puckett, Ofir Nachum, Ola Okelola,	Leclerc, Hadi Salman, Haiming Bao, Hao Sheng,	638
576	Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins,	Hart Andrin, Hessam Bagherinezhad, Hongyu Ren,	639
577	Olivier Godement, Owen Campbell-Moore, Patrick	Hunter Lightman, Hyung Won Chung, Ian Kivlichen,	640
578	Chao, Paul McMillan, Pavel Belov, Peng Su, Pe-	Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte,	641
579	ter Bak, Peter Bakkum, Peter Deng, Peter Dolan,	Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina	642
580	Peter Hoeschele, Peter Welinder, Phil Tillet, Philip	Kofman, Jakub Pachocki, James Lennon, Jason Wei,	643
581	Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming	Jean Harb, Jerry Twore, Jiacheng Feng, Jiahui Yu,	644
582	Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Ra-	Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero	645
583	jan Troll, Randall Lin, Rapha Gontijo Lopes, Raul	Candela, Joe Palermo, Joel Parish, Johannes Hei-	646
584	Puri, Reah Miyara, Reimar Leike, Renaud Gaubert,	decke, John Hallman, John Rizzo, Jonathan Gordon,	647
585	Reza Zamani, Ricky Wang, Rob Donnelly, Rob	Jonathan Uesato, Jonathan Ward, Joost Huizinga,	648
586	Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchan-	Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Ka-	649
587	dani, Romain Huet, Rory Carmichael, Rowan Zellers,	rina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood,	650
588	Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan	Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu,	651
589	Cheung, Saachi Jain, Sam Altman, Sam Schoenholz,	Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad,	652
590	Sam Toizer, Samuel Miserendino, Sandhini Agar-	Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho,	653
591	wal, Sara Culver, Scott Ethersmith, Scott Gray, Sean	Liam Fedus, Lilian Weng, Linden Li, Lindsay Mc-	654
592	Grove, Sean Metzger, Shamez Hermani, Shantanu	Callum, Lindsey Held, Lorenz Kuhn, Lukas Kon-	655
593	Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shi-	draciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd,	656
594	rong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay,	Maja Trebacz, Manas Joglekar, Mark Chen, Marko	657
595	Srinivas Narayanan, Steve Coffey, Steve Lee, Stew-	Tintor, Mason Meyer, Matt Jones, Matt Kaufner,	658
596	art Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao	Max Schwarzer, Meghan Shah, Mehmet Yatbaz,	659
597	Xu, Tarun Gogineni, Taya Christianson, Ted Sanders,	Melody Y. Guan, Mengyuan Xu, Mengyuan Yan,	660
598	Tejal Patwardhan, Thomas Cunningham, Thomas	Mia Glaese, Mianna Chen, Michael Lampe, Michael	661

662	Malek, Michele Wang, Michelle Fradin, Mike Mc-	deployment control of language models. <i>arXiv</i>	721
663	Clay, Mikhail Pavlov, Miles Wang, Mingxuan Wang,	<i>preprint arXiv:2406.15518</i> .	722
664	Mira Murati, Mo Bavarian, Mostafa Rohaninejad,		
665	Nat McAleese, Neil Chowdhury, Neil Chowdhury,	Qwen Team. 2025. Qwq-32b: Embracing the power of	723
666	Nick Ryder, Nikolas Tezak, Noam Brown, Ofir	reinforcement learning . Accessed: 2025-03-06.	724
667	Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins,		
668	Patrick Chao, Paul Ashbourne, Pavel Izmailov, Pe-	Hugo Touvron, Louis Martin, Kevin R. Stone, Peter	725
669	ter Zhokhov, Rachel Dias, Rahul Arora, Randall	Albert, Amjad Almahairi, Yasmine Babaei, Nikolay	726
670	Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Mi-	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	727
671	yara, Reimar Leike, Renny Hwang, Rhythm Garg,	Bhosale, D. Bikel, Lukas Blecher, Cristian Cantón	728
672	Robin Brown, Roshan James, Rui Shu, Ryan Cheu,	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	729
673	Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer,	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	730
674	Sam Toyer, Samuel Miserendino, Sandhini Agarwal,	Cynthia Gao, Vedanuj Goswami, Naman Goyal,	731
675	Santiago Hernandez, Sasha Baker, Scott McKinney,	A. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan,	732
676	Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani	Marcin Kardas, Viktor Kerkez, Madian Khabsa, Is-	733
677	Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang,	abel M. Kloumann, A. Korenev, Punit Singh Koura,	734
678	Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji,	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	735
679	Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	736
680	Clark, Tao Wang, Taylor Gordon, Ted Sanders, Te-	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	737
681	jal Patwardhan, Thibault Sottiaux, Thomas Degry,	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	738
682	Thomas Dimson, Tianhao Zheng, Timur Garipov,	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	739
683	Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peter-	Ruan Silva, Eric Michael Smith, R. Subramanian,	740
684	son, Tyna Eloundou, Valerie Qi, Vineet Kosaraju,	Xia Tan, Binh Tang, Ross Taylor, Adina Williams,	741
685	Vinnie Monaco, Vitchyr Pong, Vlad Fomenko,	Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan	742
686	Weiyi Zheng, Wenda Zhou, Wes McCabe, Wojciech	Zarov, Yuchen Zhang, Angela Fan, Melanie Kam-	743
687	Zaremba, Yann Dubois, Yinghai Lu, Yining Chen,	badur, Sharan Narang, Aurelien Rodriguez, Robert	744
688	Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yun-	Stojnic, Sergey Edunov, and Thomas Scialom. 2023.	745
689	yun Wang, Zheng Shao, and Zhuohan Li. 2024b.	Llama 2: Open foundation and fine-tuned chat mod-	746
690	Openai ol system card . <i>Preprint</i> , arXiv:2412.16720.	els. <i>arXiv preprint arXiv:2307.09288</i> .	747
691	Kiho Park, Yo Joong Choe, and Victor Veitch.		
692	2024. The linear representation hypothesis and	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	748
693	the geometry of large language models . <i>Preprint</i> ,	Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and	749
694	arXiv:2311.03658.	Denny Zhou. 2023. Chain-of-thought prompting elic-	750
695	Qwen, :, An Yang, Baosong Yang, Beichen Zhang,	its reasoning in large language models . <i>Preprint</i> ,	751
696	Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,	arXiv:2201.11903.	752
697	Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin,	Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017.	753
698	Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang,	Crowdsourcing multiple choice science questions .	754
699	Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang,	<i>Preprint</i> , arXiv:1707.06209.	755
700	Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li,		
701	Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji	Shu Yang, Junchao Wu, Xin Chen, Yunze Xiao, Xinyi	756
702	Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang	Yang, Derek F. Wong, and Di Wang. 2025. Under-	757
703	Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang	standing aha moments: from external observations to	758
704	Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru	internal mechanisms . <i>Preprint</i> , arXiv:2504.02956.	759
705	Zhang, and Zihan Qiu. 2025. Qwen2.5 technical		
706	report . <i>Preprint</i> , arXiv:2412.15115.	Lei Yu, Virginie Do, Karen Hambardzumyan, and	760
707	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	Nicola Cancedda. 2025a. Robust LLM safeguarding	761
708	Dario Amodei, and Ilya Sutskever. 2019. Language	via refusal feature adversarial training . In <i>The Thir-</i>	762
709	models are unsupervised multitask learners. <i>OpenAI</i>	<i>teenth International Conference on Learning Repre-</i>	763
710	<i>blog</i> .	<i>sentations</i> .	764
711	Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong,	Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,	765
712	Evan Hubinger, and Alexander Turner. 2024. Steer-	Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong	766
713	ing llama 2 via contrastive activation addition . In	Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin,	767
714	<i>Proceedings of the 62nd Annual Meeting of the As-</i>	Bole Ma, Guangming Sheng, Yuxuan Tong, Chi	768
715	<i>sociation for Computational Linguistics (Volume 1:</i>	Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jin-	769
716	<i>Long Papers)</i> , pages 15504–15522, Bangkok, Thai-	hua Zhu, Jiase Chen, Jiangjie Chen, Chengyi Wang,	770
717	land. Association for Computational Linguistics.	Hongli Yu, Weinan Dai, Yuxuan Song, Xiangpeng	771
718	Asa Cooper Stickland, Alexander Lyzhov, Jacob Pfau,	Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-	772
719	Salsabila Mahdi, and Samuel R Bowman. 2024.	Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and	773
720	Steering without side effects: Improving post-	Mingxuan Wang. 2025b. Dapo: An open-source	774
		llm reinforcement learning system at scale . <i>Preprint</i> ,	775
		arXiv:2503.14476.	776

Qingcheng Zeng, Weihao Xuan, Leyang Cui, and Rob Voigt. 2025. *Do reasoning models show better verbalized calibration?* Preprint, arXiv:2504.06564.

Anqi Zhang, Yulin Chen, Jane Pan, Chen Zhao, Aurojit Panda, Jinyang Li, and He He. 2025. *Reasoning models know when they’re right: Probing hidden states for self-verification.* Preprint, arXiv:2504.05419.

A Prompt used for Reasoning Responses Segmentation

Table 2 presents the prompt we used to query GPT-4o to segment the model’s responses into *Reasoning Segments*.

B Detailed Statistics of Segments Groups

We provide detailed statistics of the segment groups for the three datasets — GSM8K, MATH-500, and MBPP — as described in §?? in Table ?? below.

C Additional PCA Visualizations on MATH-500 and MBPP

Here, we present additional PCA results for the Self-Reflection Features on MATH-500 and MBPP in Figure 3.

D Foundations for Logit Lens Analysis

Here, we provide a brief introduction to the Logit Lens method and the background knowledge it involves.

D.1 MLP in Transformers

In transformer-based language models, the MLP is a crucial component for storing the model’s factual knowledge, and its sub-layers can be viewed as key-value memories (Geva et al., 2021). To be specific, the first layer⁵ of MLP sublayers can be viewed as a matrix W_K formed by key vectors $\{\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_n\}$, used to capture a set of patterns in the input sequence, and ultimately outputting the coefficient scores. The second layer can be viewed as a matrix W_V formed by value vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$, with each value vector containing the corresponding factual knowledge.

Formally, the output of the MLP in the transformer’s ℓ -th layer, given an input hidden state \mathbf{x}^ℓ ,

⁵In most decoder-only models, such as GPT-2 (Radford et al., 2019) and GPT-J (?), the MLP component consists of two layers, whereas in LLaMA (Touvron et al., 2023), it comprises three layers. However, we can still regard LLaMA’s first two layers collectively as the key matrices, with their output representing the coefficient scores.

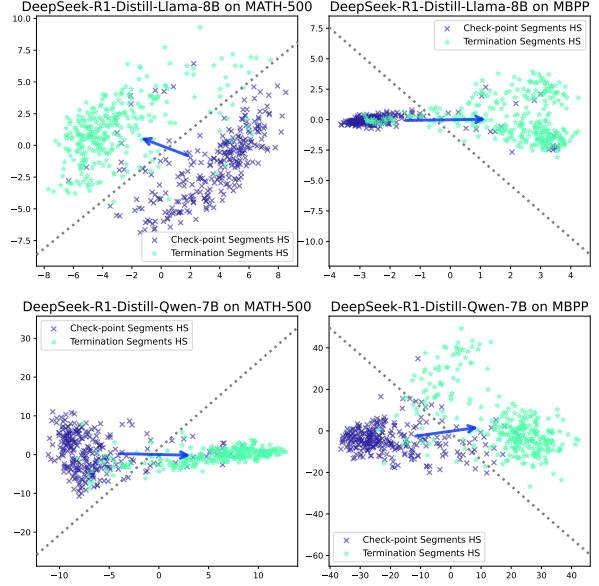


Figure 3: Visualization of the hidden states of two reasoning models on the MATH-500 and MBPP dataset using 2-dimensional PCA. The hidden states of datapoints in $\mathcal{S}_{Check-point}$ and $\mathcal{S}_{Termination}$ are positioned around the boundary (grey dashed line) fitted via logistic regression. The blue arrow approximately indicates the direction of the Self-Reflection Features. To make the image presentation clearer, we sampled 300 data points from each of $\mathcal{S}_{Check-point}$ and $\mathcal{S}_{Termination}$ for presentation. Results on other datasets are shown in §C of the Appendix.

can be defined as:

$$\mathbf{M}^\ell = f(W_K^\ell \cdot \gamma(\mathbf{x}^\ell + \mathbf{A}^\ell))W_V^\ell = \mathbf{m}^\ell W_V^\ell, \quad (3)$$

where $W_K^\ell, W_V^\ell \in \mathbb{R}^{n \times d}$. The function f and γ represent a non-linearity⁶ and layer normalization, respectively. In the transformer’s ℓ -th layer, $\mathbf{m}^\ell \in \mathbb{R}^n$ denotes the coefficient scores, and \mathbf{A}^ℓ represents the output of the attention component. The hidden state dimension is d , while the intermediate MLP has a dimension of n . Then, by denoting \mathbf{v}_j^ℓ as the j -th column (which will be called the value vector or parameter vector in the following sections) of W_V^ℓ and m_j^ℓ as the j -th element in the coefficients produced by the first layer of the MLP, we can view MLP’s output \mathbf{M}^ℓ as a linear combination of the value vectors in W_V^ℓ , with their corresponding coefficients \mathbf{m}^ℓ :

$$\mathbf{M}^\ell = \sum_{j=1}^n m_j^\ell \mathbf{v}_j^\ell, \quad (4)$$

Each m_j^ℓ here also represents the activation value of the value vector we mentioned in Table 1. Fi-

⁶For brevity, the bias term is omitted.

Prompt

- **Analyze the model response and divide into reasoning segments. Return:**

1. Labeled segments with independent solution attempts
2. Each segment must include:
 - A new full solution pathway
 - Alternative interpretations (if applicable)
 - Verification/error-checking steps (if applicable)

Format Requirements:

- Use **Segment N** headers
- Mutual exclusivity between segments
- Avoid single-step fragmentation

Examples:

Problem: “Natalia sold clips to 48 of her friends in April, and then she sold half as many clips in May. How many clips did Natalia sell altogether in April and May?”

Model Response: “[Full model response here...]”

Segmentation:

1. **Standard Calculation:**

Segment 1: Direct arithmetic approach

- April: 48 clips (given)
- May: $48/2 = 24$ clips
- Total: $48 + 24 = 72$
- Verification: $40 + 20 + 8 + 4 = 72$

2. **Algebraic Reformulation:**

Segment 2: Symbolic representation

- Let $A = 48$ (April sales)
- Define $M = A/2$ (May sales)
- Total $T = A + M = 1.5A$
- Compute $1.5 \times 48 = 72$

3. **Semantic Analysis:**

Segment 3: Ambiguity resolution

- Challenge: “sold to friends” interpretation
- Reject per-friend vs. total sales hypotheses
- Confirm $48 =$ total clips (not friends count)

Current Problem: {Problem}

Current Model Response: {Response}

Segmentation:

Table 2: Prompt for segmenting mathematical reasoning processes.

Dataset	Avg. Tokens in Termination Segments	Avg. Tokens in Check-point Segments
GSM8K	52.3	38.7
MATH-500	67.8	44.2
MBPP	41.6	29.5

Table 3: Detailed statistics of the segment groups for the three datasets — GSM8K, MATH-500, and MBPP.

nally, the hidden states at the ℓ -th layer of the language model can be defined as:

$$X^{\ell+1} = X^\ell + \mathbf{M}^\ell + \mathbf{A}^\ell, \quad (5)$$

where X^ℓ , \mathbf{M}^ℓ and \mathbf{A}^ℓ represent the hidden states, MLP’s output, and the attention component’s output in the transformer’s ℓ -th layer, respectively.

D.2 Logit Lens

[nostalgebraist \(2020\)](#); [Geva et al. \(2021\)](#) proposed that the hidden states or module parameters of a transformer-based model can be directly decoded into the vocabulary space using the model’s pre-trained unembedding matrix, enabling an investigation into the information they encode:

$$Projection = E\mathbf{v}_j^\ell, \quad (6)$$

Here, E denotes the model’s pretrained unembedding matrix, and the result of the projection, which lies in $\mathbb{R}^{|\mathcal{V}|}$, is a vector assigning a score to each token in the vocabulary \mathcal{V} . The set of the top- k highest-scoring tokens in this projection, denoted by $\mathcal{T}_{j,k}^\ell$, often reveals a clear pattern that corresponds to a specific knowledge being promoted by \mathbf{v}_j^ℓ during inference ([Geva et al., 2022b,a](#)).

E Implementation Details of the LRMs

For the inference settings of all four Large Reasoning Models, we use a temperature of 0.6, a top-p value of 0.95, and set the maximum generation length to 32,768 tokens, following the default settings.

All the experiments in this work were conducted on four 80GB NVIDIA A800 GPUs.