

# Probing the Difficulty Perception Mechanism of Large Language Models

Sunbown Lee<sup>1,2</sup>, Qingyu Yin<sup>3</sup>, Chak Tou Leong<sup>4</sup>, Jialiang Zhang<sup>1</sup>,  
Yicheng Gong<sup>2</sup>, Shiwen Ni<sup>5</sup>, Min Yang<sup>5</sup>, Xiaoyu Shen<sup>1\*</sup>

<sup>1</sup> Institute of Digital Twin, EIT <sup>2</sup> Wuhan University of Science and Technology  
<sup>3</sup> Zhejiang University <sup>4</sup> Hong Kong Polytechnic University  
<sup>5</sup> Shenzhen Institutes of Advanced Technology, CAS  
bw1863@outlook.com; xyshen@eitech.edu.cn

## Abstract

Large language models (LLMs) are increasingly deployed on complex reasoning tasks, yet little is known about their ability to internally evaluate problem difficulty, which is an essential capability for adaptive reasoning and efficient resource allocation. In this work, we investigate whether LLMs implicitly encode problem difficulty in their internal representations. Using a linear probe on the final-token representations of LLMs, we demonstrate that the difficulty level of math problems can be linearly modeled. We further locate the specific attention heads of the final Transformer layer: these attention heads have opposite activation patterns for simple and difficult problems, thus achieving perception of difficulty. Our ablation experiments prove the accuracy of the location. Crucially, our experiments provide practical support for using LLMs as automatic difficulty annotators, potentially substantially reducing reliance on costly human labeling in benchmark construction and curriculum learning. We also uncover that there is a significant difference in entropy and difficulty perception at the token level. Our study reveals that difficulty perception in LLMs is not only present but also structurally organized, offering new theoretical insights and practical directions for future research. Our code is available at <https://github.com/Aegis1863/Difficulty-Perception-of-LLMs>.

## 1 Introduction

In the context of test time scaling, a large number of peer studies are focusing on how to enable models to adaptively control the output length based on the difficulty of the problem, which can avoid lengthy reasoning on simple problems (OpenAI, 2025; Yang et al., 2025a). On the basis of DeepSeek-R1 (Guo et al., 2025), people have mastered the method of using supervised fine-tuning

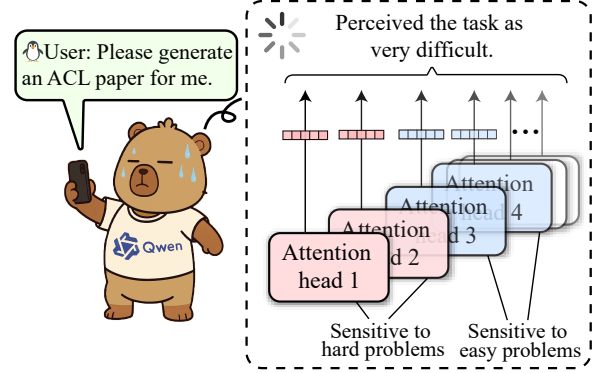


Figure 1: The LLM’s perception of problem difficulty depends on the activation state of its specific attention heads.

combined with reinforcement learning (RL) for reasoning training. Therefore, the difficulty of adaptive token budget lies in how to label difficult and simple problems, as well as how to reduce the cost of this labeling. In this study, we focus on the mechanism of the model’s perception of difficulty, as well as how to discover the simple and difficult problems that the model considers.

Evaluating the difficulty of a question is inherently complex. First, a more reliable approach often involves human annotation, which, while generally more accurate, is both expensive and highly susceptible to subjective bias, and even model and human perceptions of difficulty may differ, making consistent and scalable difficulty evaluation challenging. Moreover, many existing methods attempt to bypass human labeling by using proxy metrics, such as the length of strong large language models’ (LLMs) reasoning (Shi et al., 2025): the assumption being that longer responses indicate harder questions. However, this heuristic is frequently misleading. Powerful models like DeepSeek-R1 often produce redundant reasoning, which may reach the correct answer early but continue exploring alternative paths, inflating response length and distorting

\*Corresponding author

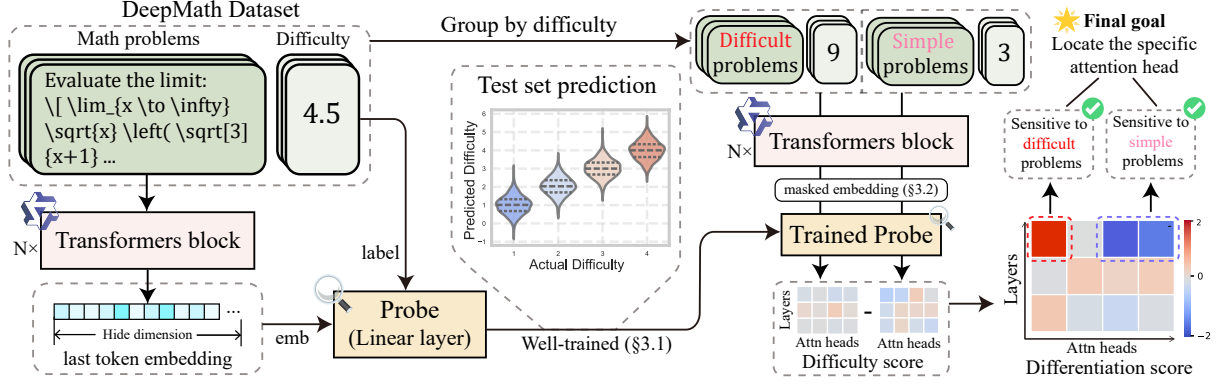


Figure 2: Probe training and attention heads pattern recognition. On the left we demonstrate how to train a difficulty probe based on last token embedding and corresponding difficulty labels. On the right, we show the process of using probe to identify attention head patterns, where the difficulty score of attention heads for difficult problems minus the difficulty score for simple problems yields a differentiation score, used to locate the attention heads most sensitive to difficulty.

difficulty estimates. This issue is especially pronounced for open-ended questions (e.g., riddles or historical), where answer length bears little relation to actual cognitive difficulty. In general, the identification of problems through existing methods is not sufficiently reliable and lacks interpretability.

To address the lack of empirical evidence for difficulty perception, we, to the best of our knowledge, are the first to explore perceiving difficulty through the internal attention heads of LLMs. Leveraging DeepMath (He et al., 2025), a high-quality mathematical benchmark annotated with difficulty levels, we empirically demonstrate that well-trained LLMs can inherently encode the difficulty of mathematical problems. In particular, some models exhibit this perception directly through distinctive attention-head activation patterns. Internal perception in LLMs can be captured via a lightweight linear probe, while identifying a question’s difficulty merely requires tracking the activations of specific attention heads. Our mechanism interpretability experiments further reveal that LLMs linearly represent mathematical difficulty in a high-dimensional embedding space, and that such representations generalize beyond the training distribution.

We conduct token-level case studies on difficulty perception, which reveals several promising research directions. The analysis explores the relationship between perceived difficulty and entropy during inference. While previous work often assumes a correlation between question difficulty and entropy, the results show that internal difficulty perception of LLMs does not consistently align with entropy variation. Tokens identified as difficult

also differ markedly from those with high entropy, suggesting that entropy-based difficulty estimation may not accurately capture the model’s internal evaluation.

In summary, the overall experimental procedure is illustrated in Figure 2 and our contributions are as follows:

- We prove that LLMs’ perception of mathematical problems is high-dimensional linear and can accurately identify this direction.
- We precisely located the difficulty perception attention head of specific LLMs. Ablation experiments show that we can modify the perception by manipulating the output of the specific attention head.
- Our case study demonstrates inconsistencies in difficulty perception and entropy, and identifies tokens that are likely to cause significant changes in difficulty perception, providing promising insights for future research.

## 2 Related Work

Our work is closely related to research on difficulty perception in efficient reasoning, which we conceptualize along two dimensions: external evaluation and internal perception.

**External evaluation.** This part of the research usually uses proxy indicators to estimate the difficulty of the problem. For example, RL-based reward functions have been used to dynamically adjust reasoning or response length according to question difficulty, balancing answer accuracy and

efficiency (Luo et al., 2025; Shen et al., 2025; Yang et al., 2025b). Similarly, Yang et al. (2025c) suggested using the shortest correct answer from a batch as a proxy for difficulty. Han et al. (2025) hypothesized that problem difficulty can be estimated via binary search on token budgets, assuming monotonic correctness with respect to token count. Several studies (Wu et al., 2025; Fang et al., 2025) trained dual reasoning modes (quick vs. slow thinking) without explicit difficulty signals, allowing the model to implicitly learn difficulty. In the aspect of suppressing overthinking, Huang et al. (2025) truncated responses once a correct answer was generated, discouraging overthinking on easy questions.

**Internal perception.** This part of the work has noticed the possible perception within the model but has not identified the attention head patterns. For example, Zhu et al. (2025b) observed clustering of easy and hard problems in the representation space of vision language models, suggesting intrinsic difficulty perception. Pioneeringly, Yin et al. (2025) used probing to monitor unsafe outputs and manipulated attention heads to improve safety, implying linearly separable safety perception. Similarly, Lee et al. (2025) found distinct representations for safe vs. dangerous prompts, indicating inherent perception capabilities. Moreover, steering has been used to guide reasoning, showing that internal representations can be effectively manipulated (Zbeeb et al., 2025; Azizi et al., 2025). In addition, exploration can be guided using token confidence, which is similar to entropy or perplexity (Ghasemabadi et al., 2025; Wang et al., 2025; Cui et al., 2025).

### 3 Method

Several peer studies have observed that, even before commencing reasoning, the embeddings generated by LLMs upon receiving a prompt can already reflect the intrinsic properties of that prompt, a phenomenon particularly pronounced in the domain of safety (Hu and Wang, 2024; Lee et al., 2025; Zheng et al., 2024). This observation has inspired our work on difficulty-perception modeling.

#### 3.1 High-dimensional Linear Probe

**Low-dimensional representation is insufficient to clearly perceive difficulty.** Although Zhu et al. (2025b) has demonstrated that multimodal models such as Qwen2.5-VL (Bai et al., 2025) and

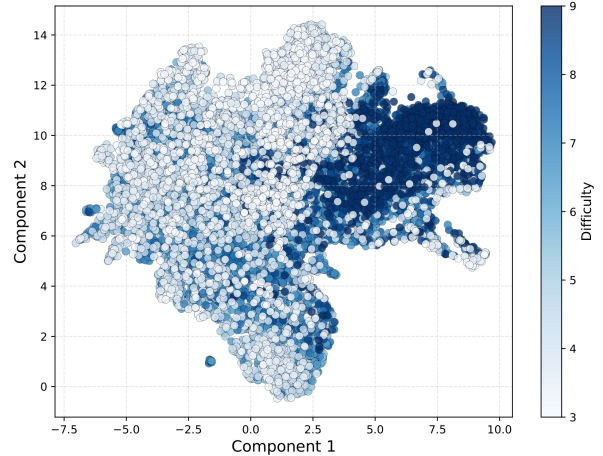


Figure 3: Qwen2.5-7B-Instruct’s low-dimensional representation for DeepMath problems and difficulty is meticulously annotated by humans. Although difficult samples seem to come together, it is generally difficult to distinguish clearly.

InternLM-VL3 (Zhu et al., 2025a) can significantly differentiate image-text pairs of different difficulty levels on low-dimensional representations at the last token embedding, this phenomenon cannot be observed intuitively in pure text scenarios. But we found that there is a difficulty-perception direction originally embedded in the high-dimensional linear space, which cannot be clearly observed in the low-dimensional space (Sheng et al., 2025), refer to Figure 3, other models’ are shown in Appendix A. However, the complex behavior of the LLMs may manifest as linear in high-dimensional embeddings. Based on the previous mechanism’s explainability foundation (Hewitt and Liang, 2019; Xu et al., 2024; Ardit et al., 2024; Tighidet et al., 2024), we believe that difficulty perception is high-dimensional linear.

**Training difficulty perception probe.** In order to investigate whether the problem embeddings encoded by the well-trained model can be consistent with the results annotated by humans, we employ a simple linear regression probe, a standard approach in representation probing (Tighidet et al., 2024). Specifically, given a  $d$ -dimensional embedding vector  $\mathbf{h} \in \mathbb{R}^d$  representing a question, the probe predicts a scalar difficulty score via a learnable linear transformation:

$$\hat{y} = \mathbf{w}^\top \mathbf{h} + b, \quad (1)$$

where  $\mathbf{w} \in \mathbb{R}^d$  and  $b \in \mathbb{R}$  are trainable parameters, and  $\hat{y}$  denotes the predicted difficulty. This probe is a standard linear regressor and is optimized by

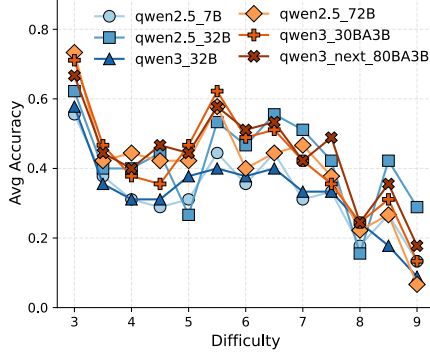


Figure 4: Accuracy rates of various models at different difficulty levels in DeepMath. As the difficulty increases, the accuracy rates of the models decrease, proving the rationality of manual labeling.

the minimum mean square error:

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2, \quad (2)$$

where  $y_i$  is ground-truth difficulty for the  $i$ -th question and  $N$  is the number of training examples. We directly adopt a regression form for modeling without introducing an activation function, in order to obtain more granular difficulty perception results.

### 3.2 Attention Head Pattern Recognition

Once the probe is well-trained, we can use it to identify specific attention head patterns.

To investigate the role of individual attention heads in encoding task difficulty, we propose a head-wise difficulty attribution framework based on selective attention ablation. Let  $\mathbf{H} \in \mathbb{R}^{B \times L \times N \times d}$  denote the multi-head attention output, where  $B$  is the batch size,  $L$  the sequence length,  $N$  the number of attention heads, and  $d$  the head dimension. The final contextual representation is obtained via the output projection  $\mathbf{W}_o \in \mathbb{R}^{(Nd) \times D}$ , yielding

$$\mathbf{Z} = \text{Reshape}(\mathbf{H})\mathbf{W}_o^\top \in \mathbb{R}^{B \times L \times D}, \quad (3)$$

with  $D = Nd$  the model’s hidden dimension.

We assume the availability of a pre-trained linear difficulty probe  $\mathbf{v}_{\text{diff}} \in \mathbb{R}^D$ , learned to predict scalar difficulty scores from the final-layer representation of the last input token. To isolate the contribution of the  $i$ -th attention head ( $i \in \{1, \dots, N\}$ ), we construct an ablated representation  $\mathbf{H}^{(i)}$  by zeroing out all heads except the

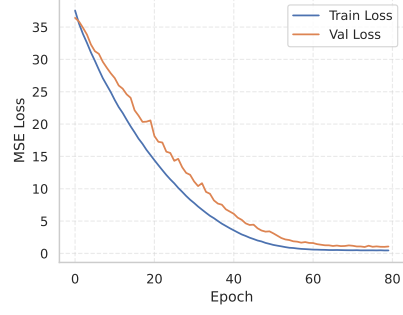


Figure 5: Probe training and validation loss. Training convergence is normal, no overfitting trend observed.

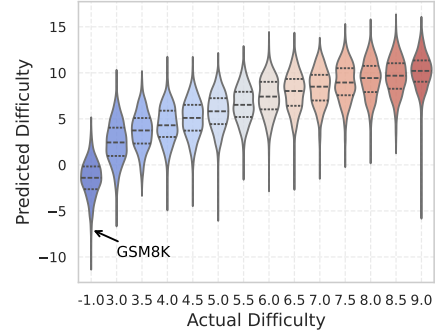


Figure 6: Probe test result. Among them, GSM8K is out of distribution data, showing lower prediction difficulty, consistent with expectations. The data on the right is the test set result of DeepMath. Since GSM8K does not have a real difficulty label, we symbolically mark it as -1 on the horizontal axis.

$i$ -th:

$$\mathbf{H}_{b,\ell,j,:}^{(i)} = \begin{cases} \mathbf{H}_{b,\ell,i,:}, & \text{if } j = i, \\ \mathbf{0}, & \text{otherwise,} \end{cases} \quad (4)$$

$$\forall b \in [B], \ell \in [L], j \in [N].$$

The corresponding projected representation is  $\mathbf{Z}^{(i)} = \text{Reshape}(\mathbf{H}^{(i)})\mathbf{W}_o^\top$ . We extract the last token embedding  $\mathbf{z}_b^{(i)} = \mathbf{z}_{b,L-1,:}^{(i)} \in \mathbb{R}^D$  for each sample  $b$ , and compute its difficulty score as the normalized projection onto the difficulty direction:

$$s_b^{(i)} = \frac{\langle \mathbf{z}_b^{(i)}, \mathbf{v}_{\text{diff}} \rangle}{\|\mathbf{v}_{\text{diff}}\|_2}. \quad (5)$$

For a batch of samples sharing the same ground-truth difficulty level, we aggregate scores across the batch to obtain the mean head-wise difficulty attribution:

$$\bar{s}^{(i)} = \frac{1}{B} \sum_{b=1}^B s_b^{(i)}. \quad (6)$$



By repeating this procedure over batches of varying difficulty (e.g., “easy” vs. “hard” instances), we derive a head-specific difficulty sensitivity profile. The discriminative power of head  $i$  is quantified by the difference in its mean attribution between high- and low-difficulty cohorts:

$$\Delta^{(i)} = \bar{s}_{\text{hard}}^{(i)} - \bar{s}_{\text{easy}}^{(i)}. \quad (7)$$

This approach enables fine-grained analysis of how individual attention heads contribute to the model’s implicit difficulty perception, revealing specialized or redundant roles across the attention subspaces.

## 4 Experiment

### 4.1 High-dimensional Linear Internal Difficulty Perception

Utilizing the DeepMath (He et al., 2025) dataset, which has meticulously annotated problem difficulty, we can train the linear probe mentioned in Subsection 3.1. We use various models (Qwen et al., 2025; Yang et al., 2025a) to test accuracy on this benchmark to verify the rationality of difficulty labels, and the experimental results are presented in Figure 4, indicating general accuracy. Another reason for using this dataset is that it was released after the Qwen2.5 release, which minimizes the risk of data leakage to the greatest extent. Training loss and verification loss reference Figure 5, exhibits a normal convergence trend without over-fitting tendencies.

We also present results on the DeepMath test set and include the GSM8K dataset (Cobbe et al., 2021) as an out-of-distribution (OOD) benchmark. The DeepMath dataset originates from mathematical olympiad problems and has been explicitly calibrated for difficulty, making it generally more challenging. In contrast, GSM8K consists of elementary school-level arithmetic problems. Although GSM8K lacks explicit difficulty annotations, its problem design suggests a substantially lower difficulty level compared to DeepMath.

The train and test results in Figure 5 and Figure 6 indicate that the probe can accurately categorize the mathematical problems in the DeepMath test set into clear difficulty levels, despite the presence of a small amount of long-tail data. At the same time, the data from GSM8K is as expected, located at a lower difficulty prediction level. We also trained a probe on llama3.1-8B-Instruct and still perform well, refer to Appendix B.1.

The above experiments on probe show that a trained LLM can linearly characterize the difficulty of the problem in high-dimensional space on last token embedding. In this case, we can obtain a difficulty-aware direction through the linear layer of the probe, and then guide the logit of the model in this direction.

### 4.2 Difficulty Perception Attention Head Localization

Since the probe can be well trained, we can further utilize the methods in Subsection 3.2 to identify attention head patterns, and refer to Figure 7 for the experimental results.

In Figure 7, we demonstrated the attention head pattern recognition results of various size models. The size of each value (color) represents the distance of the attention head pattern from difficulty level 9 to 3. Higher values (red) indicate greater sensitivity to difficult problems, while lower values (blue) indicate sensitivity to simple problems.

There are also significant differences between transformer layers. Early layers do not differ much in difficulty perception, and their pattern values change little for problems of different difficulty. However, as the depth of layers increases, their output patterns gradually become distinct.

Obviously, the pattern of the last few attention heads are more pronounced than any other layer, as it directly determines the final perception, which is also the location we need to locate. **For Qwen2.5-7B-Instruct, in the last layer, the 10th, 11th, 12th, and 13th attention heads serve to identify simple questions, while the 7th, 8th, 16th, and 23rd attention heads serve to identify difficult questions.** In practice, combination patterns of these headers work, but we can still specifically identify. To determine a specific header pattern, we can first zero out the outputs of other attention heads, then calculate the dot product between the masked embedding and the probe weights, and the difficulty perception results can be judged according to the magnitude of the numerical value.

In Figure 8, we demonstrate the differences in attention head patterns between Qwen2.5-7B-Instruct and DeepSeek-R1-Distill-Qwen-7B. We find that after DeepSeek-R1 distillation training, the pattern of the most important perception heads have been reversed, namely the 10th, 11th, 12th, and 13th attention heads, and the patterns of other heads have also undergone certain changes. In any case, the model always tends to retain the ability

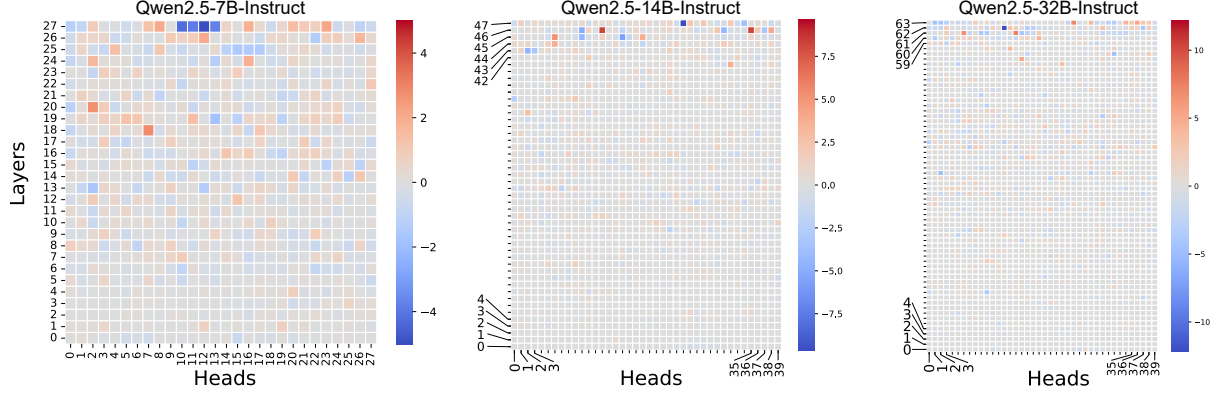


Figure 7: Different size model attention head pattern recognition results. We show the directions under significant difficulty differences (level 9 and 3). The blue attention head focuses on simple problems, while the red attention head corresponds to complex problem recognition.

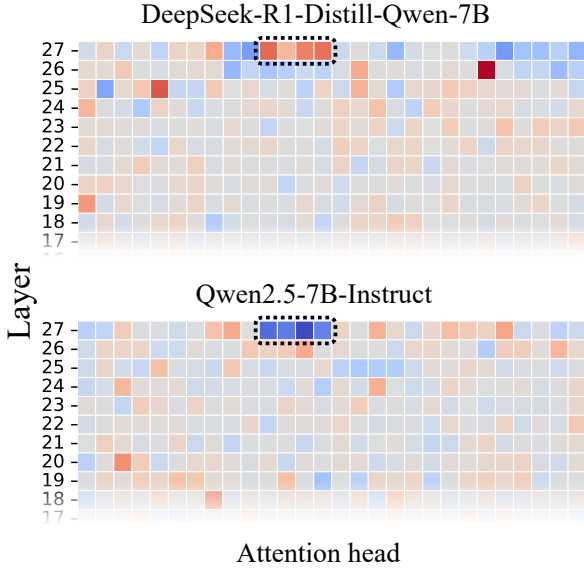


Figure 8: The attention head patterns of Qwen2.5-7B-Instruct and DeepSeek-R1-Distill-Qwen-7B. It is observed that the patterns of 4 attention heads are completely reversed, and the patterns of some other heads undergo partial changes.

to perceive difficulty and allocate this ability to different attention heads.

We also identify the attention head pattern of Llama3.1-8B-Instruct and found that its attention head is not obvious in perception of difficulty, which shows that not all models have clear perception patterns. We assume this is due to the differences in pre-training. The experimental results refer to Appendix B.2.

In the next subsection, we conduct attention head ablation experiments to verify the accuracy of the recognition results.

## 5 Attention Heads Ablation

To verify the functional specialization of attention heads revealed by our head-wise difficulty attribution analysis, we perform targeted ablation studies by intervening on heads associated with different difficulty levels. Using Qwen2.5-7B-Instruct, we isolate heads predominantly responsive to *easy* problems and those most active on *hard* ones.

We perform two complementary interventions during inference:

- **Difficulty increasing:** suppress the easy-mode heads by scaling their outputs by a factor of 0.1, while enhancing the hard-mode heads by a factor of 2.0;
- **Difficulty decreasing:** conversely, enhance the easy-mode heads ( $\times 2.0$ ) and suppress the hard-mode heads ( $\times 0.1$ ).

These manipulations are implemented by head-wise multiplication of the corresponding head outputs in the multi-head attention tensor before the output projection:

$$\mathbf{H}_{:,i,:} \leftarrow \begin{cases} \alpha_{\text{reduce}} \cdot \mathbf{H}_{:,i,:}, & i \in \mathcal{S}_{\text{easy}}, \\ \alpha_{\text{increase}} \cdot \mathbf{H}_{:,i,:}, & i \in \mathcal{S}_{\text{hard}}, \\ \mathbf{H}_{:,i,:}, & \text{otherwise,} \end{cases} \quad (8)$$

where  $\alpha_{\text{reduce}} = 0.1$ ,  $\alpha_{\text{increase}} = 2.0$ ,  $\mathcal{S}_{\text{easy}} = \{10, 11, 12, 13\}$ , and  $\mathcal{S}_{\text{hard}} = \{7, 8, 16, 23\}$ . These attention heads are determined based on the experimental results shown in Figure 7.

As shown in Table 1, the difficulty increasing setting leads to a consistent increase in the model’s estimated difficulty scores across all inputs, effectively biasing the model toward perceiving problems as more challenging. Conversely, difficulty

Real diff.	Original	Decrease	Increase
3.0	2.7	1.6 (-40.7%)	3.8 (+40.7%)
3.5	2.9	1.6 (-44.8%)	4.1 (+41.4%)
4.0	3.5	2.1 (-40.0%)	4.8 (+37.1%)
4.5	3.9	2.4 (-38.5%)	5.2 (+33.3%)
5.0	5.0	3.2 (-36.0%)	6.3 (+26.0%)
5.5	5.0	3.2 (-36.0%)	6.3 (+26.0%)
6.0	5.6	3.8 (-32.1%)	6.8 (+21.4%)
6.5	6.0	4.1 (-31.7%)	7.2 (+20.0%)
7.0	6.8	4.7 (-30.9%)	7.9 (+16.2%)
7.5	7.4	5.2 (-29.7%)	8.6 (+16.2%)
8.0	7.6	5.4 (-28.9%)	8.5 (+11.8%)
8.5	8.5	6.3 (-25.9%)	9.3 (+9.4%)
9.0	9.6	7.1 (-26.0%)	10.2 (+6.3%)

Table 1: Predicted difficulty under different adjustment conditions. Real diff. indicates the difficulty label from DeepMath’s manual annotation. Original indicates the mean difficulty estimation output by the Probe when no intervention is applied; while Decrease represents the suppression of the difficulty perception of a specific attention head, and Increase indicates the corresponding enhancement.

Real diff.	Avg. token used	
	Original	Increase
3.0	809.23	692.23 (-14.5%)
3.5	971.45	848.33 (-12.7%)
4.0	1107.70	938.10 (-15.3%)
4.5	1008.72	979.83 (-2.9%)
5.0	1059.68	1015.08 (-4.2%)

Table 2: Average token usage under different difficulties. Original indicates the state without intervention. Increase refers to the state where humans make the model believe the problem is more challenging.

decrease yields lower difficulty estimates, aligning with the hypothesis that these head groups encode complementary signals for problem complexity. These results provide causal evidence that specific attention heads are functionally specialized for perceiving inputs of different difficulty levels.

Furthermore, we observe that if we manipulate the model to perceive simple questions as more challenging, the number of output tokens decreases, as shown in Table 2. Similarly, when the original model encounters difficult problems, the number of tokens in the output will also decrease. This shows a deficit of models like Qwen2.5-7B-Instruct, that is, when facing challenging problems, they tend to give up in advance rather than make more reasoning attempts. More test results are in Appendix C.

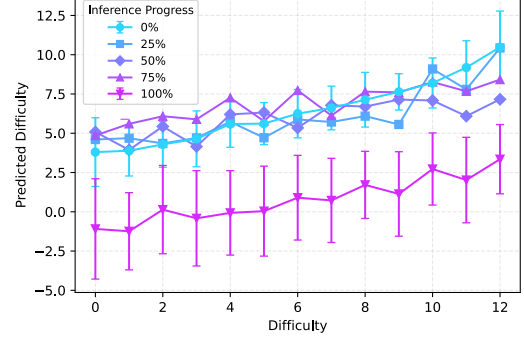


Figure 9: Utilize Probe for difficulty prediction during the inference process.

## 6 Case study

Following the release of DeepSeek-R1, a large number of research efforts have extended on the basis of GRPO (Yu et al., 2025; Zheng et al., 2025; Wang et al., 2025). Many of them are hardworking attempts to implement an adaptive token budget based on difficulty, that is, to encourage model reasoning on difficult problems and suppress model overthinking on simple problems.

In this section, we discuss the law of difficulty perception during reasoning, as well as the relationship between token-level difficulty perception and entropy. This helps us better understand how the model considers the difficulty of problems during reasoning.

### 6.1 Perception During Inference

We first collect the complete responses of Qwen2.5-7B-Instruct for samples of various difficulty math questions. Then, we retain the parts of the responses according to their length percentage as 0%, 25%, 50%, 75%, and 100%.

Referring to the experimental results in Figure 9, we demonstrated the availability and predictive patterns of the probe at various truncations during inference. It can be observed that the estimated difficulty by the probe does not change significantly during the inference process and maintained the original trend. However, at the end of inference (100%), its predicted values are significantly lower than those before, which is logical, as the problem has been solved at this point, and the model is not facing great difficulties.

### 6.2 Perceptual Difficulty and Entropy

Entropy, as an indicator to measure the determinism of the model output, is widely used to refer

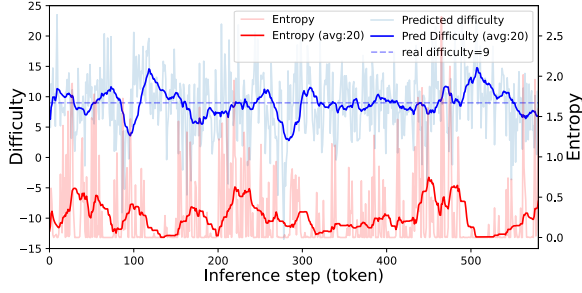


Figure 10: In the inference of a difficult math question (9.0), the change of entropy and difficulty perception.

to the difficulty of the question, which is usually true at the sentence level, but we find that it is not necessarily accurate at the token level.

In Figure 10, it is observed that the difficulty perception curve varies around the ground-truth difficulty of 9.0, while its variation characteristics are not entirely equal to the change in entropy, that is, sometimes their trends are in the same direction, and sometimes they are completely opposite. The difficulty perceived by the model are not solely determined by the entropy of the next-token probability distribution.

In Figures 11 and 12, we demonstrate an experiment of token-level difficulty perception and entropy during inference. The closer to red indicates that the model currently considers the question difficulty to be greater (Figure 11), or the uncertainty of predicting the next token to be greater (Figure 12). There is obviously a significant difference between the two, and in the perception based on difficulty, we found more granular results. An example is that our method shows greater difficulty when faced with numerical tokens, while in the results based on entropy, the entropy is mostly close to 0 when generating numbers, indicating a high degree of certainty. We believe that numbers are obviously worth paying attention to, as any error in numbers can have a huge impact on subsequent reasoning.

## 7 Conclusion

In this study, we systematically investigate how LLMs perceive problem difficulty. Based on tools and theories of mechanism interpretability, we constructed a probe to identify the high-dimensional representations of the model for problems of various difficulty levels. We trained the probe using the DeepMath dataset, which has human-labeled difficulty levels, and then used the probe to identify

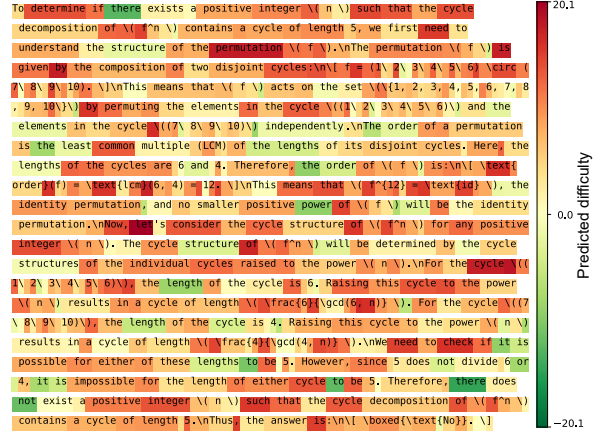


Figure 11: Token-level difficulty perception during inference.

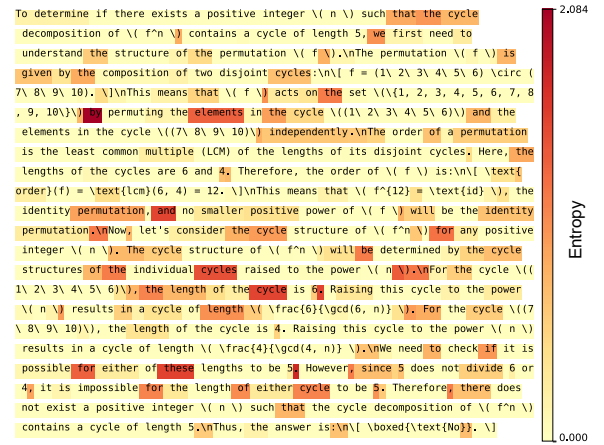


Figure 12: Token-level entropy during inference. To align with the color of Figure 11, we keep 0 as yellow.

the attention head patterns of the LLM. Ablation experiments show that our identification is accurate. Finally, we also discussed the differences between token-level difficulty perception and entropy-based methods, providing new insights for future research.

Simultaneously, we also found some promising directions. Firstly, the attention patterns of Llama3.1-8B-Instruct are very unclear compared to Qwen2.5-7B-Instruct, indicating that the quality of pre-training and post-training is likely to affect the model's ability to perceive difficulty. Secondly, the patterns of attention heads may change during training, for example, the specialized head for perceiving simple problems in DeepSeek-R1-Distill-Qwen-7B reverses compared to the original model. Finally, we found that the difficulty perception at the token level has significant differences from the characteristics of entropy, and their relationship can also be further discussed.



## Limitations

This study still has some issues worth discussing, such as whether the method based on attention head difficulty perception can be used for reward allocation in the RL training process.

## Ethics Statement and Usage Restrictions

The attention head manipulation method used in this paper is solely for the purpose of studying the internal difficulty perception of LLMs. We do not support any malicious modification or misuse of the model.

## References

- Andy Ardit, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. 2024. Refusal in language models is mediated by a single direction. *Advances in Neural Information Processing Systems*, 37:136037–136083.
- Seyedarmin Azizi, Erfan Baghaei Potraghloo, and Masoud Pedram. 2025. [Activation steering for chain-of-thought compression](#). *Preprint*, arXiv:2507.04742.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others. 2025. [Qwen2.5-vl technical report](#). *Preprint*, arXiv:2502.13923.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Yingqian Cui, Pengfei He, Jingying Zeng, Hui Liu, Xi'anfeng Tang, Zhenwei Dai, Yan Han, Chen Luo, Jing Huang, Zhen Li, Suhang Wang, Yue Xing, Jiliang Tang, and Qi He. 2025. [Stepwise perplexity-guided refinement for efficient chain-of-thought reasoning in large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 18581–18597, Vienna, Austria. Association for Computational Linguistics.
- Gongfan Fang, Xinyin Ma, and Xinchao Wang. 2025. [Thinkless: Llm learns when to think](#). *Preprint*, arXiv:2505.13379.
- Amirhosein Ghasemabadi, Keith G. Mills, Baochun Li, and Di Niu. 2025. [Guided by gut: Efficient test-time scaling with reinforced intrinsic confidence](#). *Preprint*, arXiv:2505.20325.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, and 175 others. 2025. [Deepseek-r1 incentivizes reasoning in llms through reinforcement learning](#). *Nature*, 645(8081):633–638.
- Tingxu Han, Zhenting Wang, Chunrong Fang, Shiyu Zhao, Shiqing Ma, and Zhenyu Chen. 2025. [Token-budget-aware llm reasoning](#). *Preprint*, arXiv:2412.18547.
- Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu, Zhenwen Liang, Wenxuan Wang, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2025. [Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical dataset for advancing reasoning](#). *Preprint*, arXiv:2504.11456.
- John Hewitt and Percy Liang. 2019. [Designing and interpreting probes with control tasks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.
- Leyang Hu and Boran Wang. 2024. [Droj: A prompt-driven attack against large language models](#). *Preprint*, arXiv:2411.09125.
- Chengyu Huang, Zhengxin Zhang, and Claire Cardie. 2025. [Hapo: Training language models to reason concisely via history-aware policy optimization](#). *Preprint*, arXiv:2505.11225.
- Sunbowen Lee, Shiwen Ni, Chi Wei, Shuaimin Li, Liyang Fan, Ahmadreza Argha, Hamid Alinejad-Rokny, Ruifeng Xu, Yicheng Gong, and Min Yang. 2025. [xjailbreak: Representation space guided reinforcement learning for interpretable llm jailbreaking](#). *Preprint*, arXiv:2501.16727.
- Haotian Luo, Li Shen, Haiying He, Yibo Wang, Shiwei Liu, Wei Li, Naiqiang Tan, Xiaochun Cao, and Dacheng Tao. 2025. [O1-pruner: Length-harmonizing fine-tuning for o1-like reasoning pruning](#). *Preprint*, arXiv:2501.12570.
- OpenAI. 2025. [gpt-oss-120b & gpt-oss-20b model card](#). *Preprint*, arXiv:2508.10925.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, and 25 others. 2025. [Qwen2.5 technical report](#). *Preprint*, arXiv:2412.15115.
- Yi Shen, Jian Zhang, Jieyun Huang, Shuming Shi, Wenjing Zhang, Jiangze Yan, Ning Wang, Kai Wang, Zhaoxiang Liu, and Shiguo Lian. 2025. [Dast: Difficulty-adaptive slow-thinking for large reasoning models](#). *Preprint*, arXiv:2503.04472.

- Leheng Sheng, An Zhang, Zijian Wu, Weixiang Zhao, Changshuo Shen, Yi Zhang, Xiang Wang, and Tat-Seng Chua. 2025. On reasoning strength planning in large reasoning models. *CoRR*, abs/2506.08390.
- Taiwei Shi, Yiyang Wu, Linxin Song, Tianyi Zhou, and Jieyu Zhao. 2025. [Efficient reinforcement fine-tuning via adaptive curriculum learning](#). *Preprint*, arXiv:2504.05520.
- Parshin Shojaei, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. 2025. [The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity](#). *Preprint*, arXiv:2506.06941.
- Zineddine Tighidet, Jiali Mei, Benjamin Piwowarski, and Patrick Gallinari. 2024. [Probing language models on their knowledge source](#). In *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 604–614, Miami, Florida, US. Association for Computational Linguistics.
- Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin Yang, Zhenru Zhang, Yuqiong Liu, An Yang, Andrew Zhao, Yang Yue, Shiji Song, Bowen Yu, Gao Huang, and Junyang Lin. 2025. [Beyond the 80/20 rule: High-entropy minority tokens drive effective reinforcement learning for llm reasoning](#). *Preprint*, arXiv:2506.01939.
- Siye Wu, Jian Xie, Yikai Zhang, Aili Chen, Kai Zhang, Yu Su, and Yanghua Xiao. 2025. [Arm: Adaptive reasoning model](#). *Preprint*, arXiv:2505.20258.
- Heming Xia, Chak Tou Leong, Wenjie Wang, Yongqi Li, and Wenjie Li. 2025. [Tokenskip: Controllable chain-of-thought compression in llms](#). *Preprint*, arXiv:2502.12067.
- Zhihao Xu, Ruixuan Huang, Changyu Chen, and Xiting Wang. 2024. [Uncovering safety risks of large language models through concept activation vector](#). *Preprint*, arXiv:2404.12038.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025a. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- Junjie Yang, Ke Lin, and Xing Yu. 2025b. [Think when you need: Self-adaptive chain-of-thought learning](#). *Preprint*, arXiv:2504.03234.
- Wenkai Yang, Shuming Ma, Yankai Lin, and Furu Wei. 2025c. [Towards thinking-optimal scaling of test-time compute for llm reasoning](#). *Preprint*, arXiv:2502.18080.
- Qingyu Yin, Chak Tou Leong, Linyi Yang, Wenxuan Huang, Wenjie Li, Xiting Wang, Jaehong Yoon, YunXing, XingYu, and Jinjin Gu. 2025. [Refusal falls off a cliff: How safety alignment fails in reasoning?](#) *Preprint*, arXiv:2510.06036.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, and 16 others. 2025. [Dapo: An open-source llm reinforcement learning system at scale](#). *Preprint*, arXiv:2503.14476.
- Mohammad Zbeeb, Hasan Abed Al Kader Hammoud, and Bernard Ghanem. 2025. [Reasoning vectors: Transferring chain-of-thought capabilities via task arithmetic](#). *Preprint*, arXiv:2509.01363.
- Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, Jingren Zhou, and Junyang Lin. 2025. [Group sequence policy optimization](#). *Preprint*, arXiv:2507.18071.
- Chujie Zheng, Fan Yin, Hao Zhou, Fandong Meng, Jie Zhou, Kai-Wei Chang, Minlie Huang, and Nanyun Peng. 2024. [On prompt-driven safeguarding for large language models](#). *Preprint*, arXiv:2401.18018.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, Zhangwei Gao, Erfei Cui, Xuehui Wang, Yue Cao, Yangzhou Liu, Xingguang Wei, Hongjie Zhang, Haomin Wang, Weiye Xu, and 32 others. 2025a. [InternV13: Exploring advanced training and test-time recipes for open-source multimodal models](#). *Preprint*, arXiv:2504.10479.
- Yubo Zhu, Dongrui Liu, Zecheng Lin, Wei Tong, Sheng Zhong, and Jing Shao. 2025b. [The llm already knows: Estimating llm-perceived question difficulty via hidden representations](#). *Preprint*, arXiv:2509.12886.

## A Low-dimensional Representation

We also demonstrated the low-dimensional representations of the DeepSeek-R1-Distill-Qwen-7B and TokenSkip (Xia et al., 2025), finding it difficult to clearly distinguish the spatial locations of different difficulty of questions. TokenSkip is a method of simplifying output by pruning tokens, allowing the model to respond to questions with more concise language. We tested TokenSkip to observe its impact on the model’s difficulty perception. The experimental results show that there is little effect.

## B Experiment Results of Llama3.1

We performed the same experiment on Llama3.1-8B-Instruct, which is shown in this section.

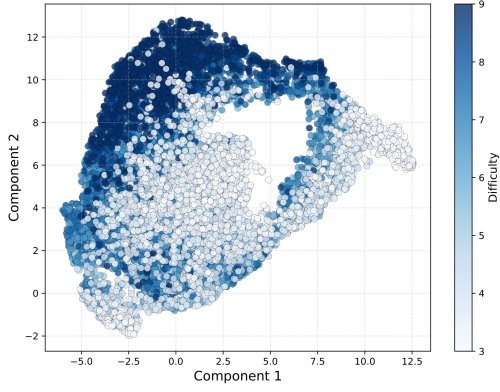


Figure 13: DeepSeek-R1-Distill-Qwen-7B’s low-dimensional representation for DeepMath problems and difficulty is meticulously annotated by humans.

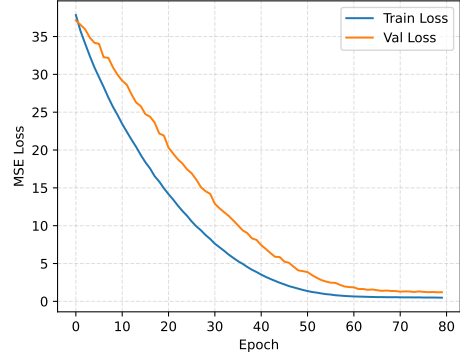


Figure 15: Llama3.1-8B-Instruct’s probe training and validation loss.

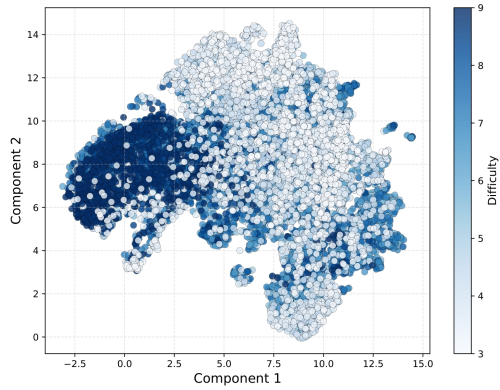


Figure 14: Tokenskip’s low-dimensional representation for DeepMath problems and difficulty is meticulously annotated by humans.

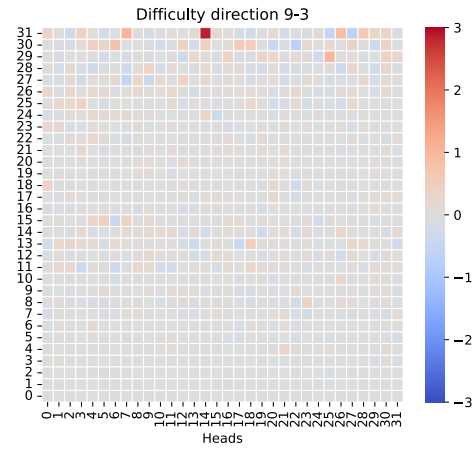


Figure 16: Llama3.1-8B-Instruct’s attention head pattern recognition.

## B.1 Probe Training

Llama3.1-8B-Instruct’s probe training still has an good effect, refer to Figure 15.

## B.2 Attention Head Pattern

The Llama’s attention head pattern is very insignificant, as referenced in Figure 16. We can only observe that the 14th attention head in the last layer obviously serves to identify difficult problems. This is different from the Qwen series model, and we assume the reason is the pre-training quality.

## C Reasoning Collapse Token Reduction

Referring to Figure 16, we found the reasoning collapse token reduction phenomenon. This means that when some models encounter very difficult problems, the number of tokens output dropped significantly, which is very similar to the findings of Shojaei et al. (2025).

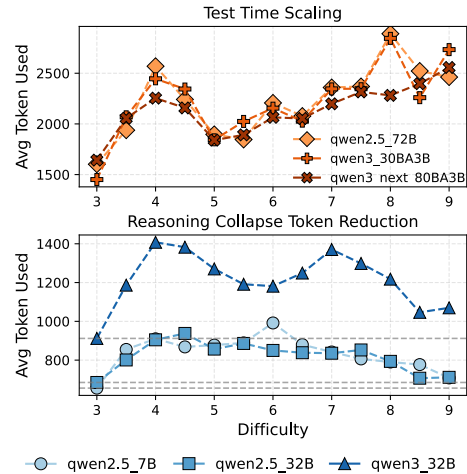


Figure 17: Reasoning Collapse Token Reduction. As the problems faced by the model become more difficult, the number of tokens it uses gradually increases. However, when some models reach a certain level of difficulty, the reasoning tokens they spend significantly decreased, while the other models continued to increase.