

# What Affects the Effective Depth of Large Language Models?

Anonymous ACL submission

## Abstract

The scaling of large language models (LLMs) emphasizes increasing depth, yet performance gains diminish with added layers. Prior work introduces the concept of “effective depth”, arguing that deeper models fail to fully utilize their layers for meaningful computation. Building on this, we systematically study how effective depth varies with model scale, training type, and task difficulty. First, we analyze the model behavior of Qwen-2.5 family (1.5B-32B) and find that while the number of effective layers grows with model size, the effective depth ratio remains stable. Besides, comparisons between base and corresponding long-CoT models show no increase in effective depth, suggesting that improved reasoning stems from longer context rather than deeper per-token computation. Furthermore, evaluations across tasks of varying difficulty indicate that models do not dynamically use more layers for harder problems. Our results suggest that current LLMs underuse available depth across scales, training paradigms and tasks of varying difficulties, pointing out research opportunities on increasing the layer utilization rate of LLMs, model pruning, and early exiting. Our code is released at [https://anonymous.4open.science/r/what\\_affects\\_effective\\_depth-4349](https://anonymous.4open.science/r/what_affects_effective_depth-4349).

## 1 Introduction

The scaling of large language models (LLMs) (Yang et al., 2025; Qwen Team, 2025; Grattafiori et al., 2024; DeepSeek-AI, 2024; Achiam et al., 2023) has consistently emphasized increased depth, with empirical evidence suggesting that model performance improves with additional layers—despite diminishing returns. As pointed out by Csordás et al. (2025), this trend raises a fundamental question: are these models truly leveraging their depth to perform more complex, hierarchical computations, or are they merely distributing similar computational operations over a greater number of layers?

Csordás et al. (2025) reveals a striking underutilization of depth from a mechanistic perspective: layers in the second half are simply refining existing representations rather than contributing to novel feature composition or conducting deeper reasoning. The study introduces the concept of “effective depth” and suggests that inefficient depth utilization may be a fundamental cause of diminishing scaling returns. Building directly upon this foundation, our work seeks to systematically investigate the factors that influence this effective depth. We aim to achieve a more comprehensive understanding of how depth utilization behaves across model scale, specialized training, and task difficulty. Our findings are as follows:

**Regarding model size.** Following the methodologies established in prior work, we first analyze the Qwen-2.5 model family (from 1.5B to 32B) (Qwen Team, 2025) using a suite of techniques including residual cosine similarity, logit lens, layer effects on future computations, residual erasure and integrated gradients (Csordás et al., 2025; Nostalgebraist, 2020). Our results confirm the core phenomenon: there exists a phase transition where early layers drive feature composition and later layers engage in minor refinements. Furthermore, while the absolute number of these “effective” layers increases with model size, the ratio of effective depth to total depth remains stable. This aligns with the conclusions of Csordás et al. (2025) that larger models do not fundamentally alter their computational strategy; they simply replicate the same utilization pattern over a larger number of layers, rather than using the extra depth to invent new types of computation. This finding provides a nuanced explanation for diminishing returns: wider models gain new capabilities, while deeper models primarily gain precision.

**Regarding long-CoT models.** Given that long-CoT models have demonstrated exceptional per-

formance in complex reasoning tasks (DeepSeek-AI, 2025; OpenAI, 2024), a natural hypothesis is that they might achieve this by more effectively exploiting their depth for “deeper” reasoning in each forward pass. To test this, we compare the effective depth of base and instruct models in the Qwen-2.5 model family (Qwen Team, 2025) against their corresponding DeepSeek-R1-distill counterparts (DeepSeek-AI, 2025). Surprisingly, our analysis reveals no significant increase in effective depth. The enhanced reasoning performance appears not to be driven by a fundamental change in how the model utilizes its layers during each forward pass. Instead, the gains are likely attributable to the model’s optimized ability to reason over longer sequences, not to deeper computation within a single token’s forward process.

**Regarding task difficulty.** We further probe whether models dynamically allocate their depth based on computational demand. One might expect harder problems to require and therefore activate deeper layers. We evaluate models on a difficulty spectrum from HellaSwag (natural language understanding) (Zellers et al., 2019) to GSM8K (grade school math) (Cobbe et al., 2021) to AIME24 (high school math contests) (MAA). Counter-intuitively, the effective depth remains largely consistent across all tasks, regardless of their varying difficulty.

In summary, current LLMs, across scales, specialized training regimes and task difficulties, fail to fully exploit their available depth.

## 2 Preliminary

We mainly focus on the Qwen-2.5 model family (Qwen Team, 2025) (including base models and instruct models), and their corresponding DeepSeek-R1-Distill versions (DeepSeek-AI, 2025). They are all pre-norm Transformers (Xiong et al., 2020; Vaswani et al., 2017) and the forward process of a layer  $l$  is as follows:

$$\mathbf{a}_l = \text{SelfAttention}_l(\text{RMSNorm}(\mathbf{h}_l)) \quad (1)$$

$$\hat{\mathbf{h}}_l = \mathbf{h}_l + \mathbf{a}_l \quad (2)$$

$$\mathbf{m}_l = \text{MLP}_l(\text{RMSNorm}(\hat{\mathbf{h}}_l)) \quad (3)$$

$$\mathbf{h}_{l+1} = \hat{\mathbf{h}}_l + \mathbf{m}_l \quad (4)$$

Here,  $\mathbf{h}_l \in \mathbb{R}^{n_{\text{context}} \times d_{\text{model}}}$  is the residual stream (Elhage et al., 2021),  $\mathbf{a}_l, \mathbf{m}_l$  are the outputs of the SelfAttention layers and MLP layers,

which are directly added back to the residual stream.  $n_{\text{context}}$  is the length of the input sequence, and  $d_{\text{model}}$  is the dimension of the hidden states of the model. RMSNorm (Zhang and Sennrich, 2019) is adopted in the Qwen-2.5 model family to replace traditional layer normalization (Xiong et al., 2020). Following Csordás et al. (2025), we denote  $\text{SelfAttention}_l(\cdot)$  and  $\text{MLP}_l(\cdot)$  as “sublayers”.

The residual stream starts with  $\mathbf{h}_0 = \text{Embedding}(x)$ , where  $x \in \mathbb{N}^{n_{\text{context}}}$  is the sequence of token\_ids. It then passes through the output layer, producing the output probability distribution over vocabulary:  $\mathbf{y} = \text{softmax}(\text{RMSNorm}(\mathbf{h}_L)\mathbf{W}^{\text{out}})$ , where  $\mathbf{y} \in \mathbb{R}^{n_{\text{context}} \times |V|}$ ,  $\mathbf{W}^{\text{out}} \in \mathbb{R}^{d_{\text{model}} \times |V|}$ ,  $L$  is the number of layers in the model,  $V$  is the vocabulary.

## 3 Methods

Csordás et al. (2025) proposes a suite of methods to qualitatively probe effective depth. We introduce and extend these methods to qualitatively assess effective depth across different models and datasets:

**Residual cosine similarity.** Residual cosine similarity measures how each layer or sublayer interacts with the residual stream. For a given layer  $l$ , we compute the cosine similarity between its contribution (the output of either SelfAttention  $\mathbf{a}_l$ , MLP  $\mathbf{m}_l$ , or their sum) and the resulting residual state  $\mathbf{h}_l$ . Formally, the similarities are defined as  $\text{cosim}(\mathbf{a}_l + \mathbf{m}_l, \mathbf{h}_l)$  for the full layer,  $\text{cosim}(\mathbf{a}_l, \mathbf{h}_l)$  for self-attention, and  $\text{cosim}(\mathbf{m}_l, \mathbf{h}_l + \mathbf{a}_l)$  for the MLP. The intuition is that a cosine similarity near zero suggests the module writes a new, orthogonal feature into the residual stream; negative values indicate feature erasure; and positive values signify the amplification of an existing feature.

**Logit Lens.** Logit lens evaluates how early the models output distribution begins to stabilize. We decode the hidden state  $\mathbf{h}_l$  using the models output projection and compute the KL divergence between this early distribution and the models final distribution. Additionally, we measure the overlap between the top-5 tokens from this intermediate distribution and from the final distribution.

**Layer effects on future computation.** Here we probe the influence of skipping a layer on subsequent computations. For a given prompt, we first run a forward pass to record the residual states  $\mathbf{h}_l$ . We then intervene by skipping a specific layer  $s$  for all token positions  $t \leq t_s$  (where  $t_s$  is a sampled position within the sequence), effectively

setting  $\bar{h}_{s+1} := \bar{h}_s$  for those tokens. The effect of this intervention is measured on the subsequent tokens ( $t > t_s$ ) by computing the relative change in the contribution of a later layer  $l > s$ :  $\|(h_{l+1} - h_l) - (\bar{h}_{l+1} - \bar{h}_l)\|_2 / \|h_{l+1} - h_l\|_2$ . The maximum value of this metric across multiple prompts and sequence positions is taken. We also compare the final output probabilities via  $\|y - \bar{y}\|_2$ .

**Residual erasure.** Residual erasure identifies until which layer information from a specific token remains relevant for the final prediction. For a token at position  $t$  and layer  $l$ , we intervene by replacing its residual vector  $h_{l+1}[t]$  with an uninformative baseline the average residual vector at layer  $l$  computed over a dataset (GSM8K here), while leaving all other tokens unchanged. The effect is quantified as the maximum change in prediction norm ( $\|y - \bar{y}\|_2$ ) among all answer tokens.

**Integrated gradients.** The metric attributes the models prediction on the answer tokens to contributions from each layer. We compute the gradients of the output logits for all answer tokens to the activation at each layer and each token position.

Beyond these qualitative probes, we introduce two quantitative measures to compare effective depth across models and datasets. For *residual cosine similarity*, we average the similarity scores across layers, MLPs, and SelfAttention modules, and identify the effective depth as the point where the averaged similarity transitions from negative to positive. For the *logit lens*, we use two metrics: we define the effective depth as the layer where the KL divergence from the final output drops below half of its maximum observed value, and alternatively, as the layer where the top-5 token overlap with the final output first exceeds 0.3.

## 4 Experiments

Main experiments presented are performed with Qwen-2.5 model family (Qwen Team, 2025) and their corresponding reasoning models provided by DeepSeek-AI (2025), using NDIF and NNsight (Fiotto-Kaufman et al., 2024). Following Csordás et al. (2025), for all the methods, unless specified otherwise, the results are computed on 10 random examples from specified datasets.

### 4.1 Does Model Size Affect Effective Depth?

The residual cosine similarity, shown in Figure 1, exhibits a consistent pattern across models: an initial positive phase declines into negative values

before returning to positive. The initial near-zero similarity in shallow layers suggests context integration, while the subsequent positive phase corresponds to feature refinement. The first half of the network is predominantly characterized by feature erasure (negative similarity), until a sharp phase transition occurs near the middle layers, after which the model begins strengthening existing features.

We quantify the corresponding depth of this transition in Table 1 (Cosine Similarity). The results show that the effective depth ratio remains remarkably stable. This indicates that larger models contain a growing number of “ineffective” layers that do not contribute to feature composition.

The logit lens analysis, as shown in Figure 2, further supports this conclusion. The KL divergence between intermediate and final predictions shows a sharp drop in the second half of the network, while the top-5 token overlap exhibits a concurrent sharp rise. Together, these indicate a transition from computation to refinement. As quantified in Table 1, the depth of this transition, measured both by KL divergence (half-max point) and overlap (exceeding 0.3), is slightly less consistent across scales than the cosine similarity metric, with a mild increasing trend in ratio for larger models.

Furthermore, the effect of skipping layers on downstream computations, illustrated in Figure 3, reveals that layers in the second half have substantially less influence on both later layers and final output predictions. This pattern is consistent across all model sizes, with similar decay profiles.

Finally, results from integrated gradients (Figure 4(a)) and residual erasure (Figure 4(b)) show that the dependence of answer token predictions on earlier layers declines markedly in the second half of the network. The position of this decline remains stable relative to network depth across model sizes.

### 4.2 Do Long-CoT Models Think Deeper?

Given that long-CoT models demonstrate superior performance on complex reasoning (DeepSeek-AI, 2025; OpenAI, 2024), one might hypothesize that they achieve this by utilizing deeper computations within each forward pass. To test this, we compare the effective depth of DeepSeek-R1-Distill models (DeepSeek-AI, 2025) against their corresponding base models (Qwen Team, 2025). As summarized in Table 1, we find no significant difference in effective depth ratio between long-CoT and base models. This consistency is further illustrated across all probing methods: residual cosine

Table 1: Effective depth (ED) and effective depth ratio (ratio =  $\frac{ED+1}{L}$ ) across base, instruct, and long-CoT models of different sizes (1.5B to 32B parameters) and on datasets with varying difficulty.

	Cosine Similarity						Logit Lens KL						Logit Lens Overlap					
	HellaSwag		GSM8K		AIME24		HellaSwag		GSM8K		AIME24		HellaSwag		GSM8K		AIME24	
	ED	ratio	ED	ratio	ED	ratio	ED	ratio	ED	ratio	ED	ratio	ED	ratio	ED	ratio	ED	ratio
DS-R1-Qwen-1.5B	17	0.64	16	0.61	17	0.64	20	0.75	1	0.07	24	0.89	23	0.86	23	0.86	24	0.89
Qwen2.5-1.5B-Instruct	16	0.61	20	0.75	19	0.71	21	0.79	22	0.82	23	0.86	23	0.86	23	0.86	23	0.86
Qwen2.5-Math-1.5B	16	0.61	16	0.61	16	0.61	20	0.75	22	0.82	23	0.86	23	0.86	23	0.86	23	0.86
DS-R1-Qwen-7B	16	0.61	16	0.61	16	0.61	24	0.89	24	0.89	24	0.89	25	0.93	25	0.93	24	0.89
Qwen2.5-7B-Instruct	17	0.64	20	0.75	18	0.68	25	0.93	25	0.93	25	0.93	26	0.96	26	0.96	26	0.96
Qwen2.5-Math-7B	16	0.61	11	0.43	16	0.61	23	0.86	23	0.86	23	0.86	24	0.89	24	0.89	24	0.89
DS-R1-Qwen-14B	26	0.56	30	0.65	30	0.65	40	0.85	39	0.83	41	0.88	44	0.94	44	0.94	44	0.94
Qwen2.5-14B-Instruct	27	0.58	32	0.69	30	0.65	40	0.85	41	0.88	42	0.90	45	0.96	45	0.96	45	0.96
Qwen2.5-14B	27	0.58	30	0.65	30	0.65	40	0.85	40	0.85	42	0.90	45	0.96	45	0.96	45	0.96
DS-R1-Qwen-32B	42	0.67	42	0.67	46	0.73	58	0.92	55	0.88	57	0.91	61	0.97	58	0.92	58	0.92
Qwen2.5-32B-Instruct	43	0.69	46	0.73	43	0.69	60	0.95	58	0.92	58	0.92	61	0.97	60	0.95	60	0.95
Qwen2.5-32B	43	0.69	46	0.73	46	0.73	60	0.95	57	0.91	59	0.94	61	0.97	59	0.94	60	0.95

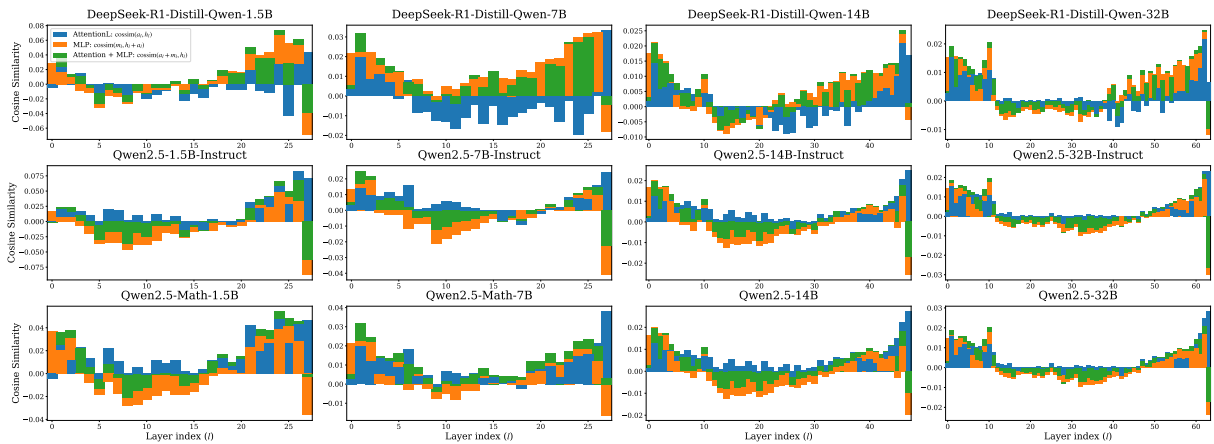


Figure 1: Cosine similarity of (sub)layer contributions and the residual evaluated on GSM8K.

similarity (Figure 1), logit lens (Figure 2), layer-skipping effects (Figure 3), integrated gradients (Figure 4(a)), and residual erasure (Figure 4(b)). The results are consistent that long-CoT models do not exhibit a deeper utilization of the network. Instead, the performance gains appear to stem from the models enhanced ability to reason over longer sequences, suggesting that “wider” context, rather than “deeper” per-token computation, underlies these improvements.

### 4.3 Does Task Difficulty Affect Effective Depth?

We next investigate whether models dynamically adjust their effective depth in response to computational demand, expecting that harder tasks might engage deeper layers. We evaluate models on three tasks of increasing difficulty: HellaSwag (natural language understanding) (Zellers et al., 2019), GSM8K (grade school math) (Cobbe et al., 2021), and AIME24 (high school math contests) (MAA). Results in Table 1 show that effective depth re-

mains largely consistent across all tasks, indicating that model depth utilization is not adaptive to problem difficulty. Additional results are provided in Appendix C, including residual cosine similarity (Figure 5), effects of skipping layers on future computations (Figure 6) and output distributions (Figure 7), as well as logit lens KL divergence (Figure 8) and token overlap (Figure 9).

## 5 Conclusion

Our study provides a comprehensive analysis of the factors influencing effective depth in LLMs, examining model scale, training strategies, and task difficulty. First, the ratio of effective depth remains stable as model size increases. Second, while long-CoT training enhances reasoning performance, it does not lead to an increase in effective depth. Third, effective depth remains consistent across tasks of varying difficulty, suggesting that models do not dynamically allocate computational depth based on problem complexity.

## 322 Limitations

323 This work follows the methodology of Csordás  
324 et al. (2025) to comprehensively analyze factors  
325 influencing effective depth. We introduce quantita-  
326 tive metrics based on residual cosine similarity and  
327 logit lens to compare effective depth across models  
328 and datasets. However, the proposed metrics—  
329 particularly the two variants of logit lens—remain  
330 relatively straightforward and exhibit some insta-  
331 bility. Developing more robust and well-validated  
332 measures of effective depth is an important direc-  
333 tion for future research.

334 Furthermore, while we confirm and extensively  
335 analyze the phenomenon of depth under-utilization  
336 across model scales, training strategies, and task  
337 demands, this study does not propose solutions to  
338 improve layer utilization. Our findings highlight  
339 the need for future work to explore architectural  
340 or training approaches that enable models to lever-  
341 age their full depth more effectively. Initial efforts  
342 along these lines have emerged in the direction of  
343 modifying model architectures (Sun et al., 2025; Li  
344 et al., 2024; Kapl et al., 2025), yet further investi-  
345 gation is warranted.

## 346 Ethical Considerations

347 This work presents a diagnostic analysis of the in-  
348 ternal computational patterns in LLMs. Our study  
349 is based entirely on publicly available, open-source  
350 models (the Qwen-2.5 family and DeepSeek-R1-  
351 distill model family) and standard, open bench-  
352 marks. As such, this research does not involve the  
353 collection of new data, the creation of new models,  
354 or any direct deployment. All code and analysis  
355 methods are released to ensure reproducibility and  
356 transparency.

357 AI assistants were utilized for language polish-  
358 ing and refinement, strictly limited to improving the  
359 fluency and clarity the text. All technical content,  
360 experimental results, analyses, and conclusions re-  
361 main the original work of the authors.

## 362 References

363 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
364 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
365 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
366 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-  
367 cal report. *arXiv preprint arXiv:2303.08774*.

368 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
369 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
370 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro

Nakano, and 1 others. 2021. Training verifiers  
371 to solve math word problems. *arXiv preprint*  
372 *arXiv:2110.14168*. 373

Róbert Csordás, Christopher D Manning, and Christo-  
374 pher Potts. 2025. Do language models use their depth  
375 efficiently? *arXiv preprint arXiv:2505.13898*. 376

DeepSeek-AI. 2024. *Deepseek-v3 technical report*.  
377 *Preprint*, arXiv:2412.19437. 378

DeepSeek-AI. 2025. *Deepseek-r1: Incentivizing rea-  
379 soning capability in llms via reinforcement learning*.  
380 *Preprint*, arXiv:2501.12948. 381

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom  
382 Henighan, Nicholas Joseph, Ben Mann, Amanda  
383 Askell, Yuntao Bai, Anna Chen, Tom Conerly, and  
384 1 others. 2021. A mathematical framework for  
385 transformer circuits. *Transformer Circuits Thread*,  
386 1(1):12. 387

Jaden Fiotto-Kaufman, Alexander R Loftus, Eric Todd,  
388 Jannik Brinkmann, Caden Juang, Koyena Pal, Can  
389 Rager, Aaron Mueller, Samuel Marks, Arnab Sen  
390 Sharma, and 1 others. 2024. Nnsight and ndif: De-  
391 mocratizing access to foundation model internals.  
392 *arXiv preprint arXiv:2407.14561*. 393

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,  
394 Abhinav Pandey, Abhishek Kadian, Ahmad Al-  
395 Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,  
396 Alex Vaughan, and 1 others. 2024. *The llama 3 herd  
397 of models*. *Preprint*, arXiv:2407.21783. 398

Ferdinand Kapl, Emmanouil Angelis, Tobias Höppe,  
399 Kaitlin Maile, Johannes von Oswald, Nino Scherrer,  
400 and Stefan Bauer. 2025. Do depth-grown models  
401 overcome the curse of depth? an in-depth analysis.  
402 *arXiv preprint arXiv:2512.08819*. 403

Pengxiang Li, Lu Yin, and Shiwei Liu. 2024. Mix-In:  
404 Unleashing the power of deeper layers by combining  
405 pre-In and post-In. *arXiv preprint arXiv:2412.13795*. 406

MAA. *American invitational mathematics examination-  
407 aime 2024, 2024*. 408

Nostalgebraist. 2020. *Interpreting gpt: The logit lens*. 409

OpenAI. 2024. *Learning to reason with llms*. 410

Qwen Team. 2025. *Qwen2.5 technical report*. *Preprint*,  
411 arXiv:2412.15115. 412

Wenfang Sun, Xinyuan Song, Pengxiang Li, Lu Yin,  
413 Yefeng Zheng, and Shiwei Liu. 2025. The curse  
414 of depth in large language models. *arXiv preprint*  
415 *arXiv:2502.05795*. 416

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob  
417 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
418 Kaiser, and Illia Polosukhin. 2017. Attention is all  
419 you need. *Advances in neural information processing*  
420 *systems*, 30. 421

422	Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng,
423	Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan
424	Lan, Liwei Wang, and Tiejian Liu. 2020. On layer
425	normalization in the transformer architecture. In <i>In-</i>
426	<i>ternational conference on machine learning</i> , pages
427	10524–10533. PMLR.
428	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,
429	Binyuan Hui, Bo Zheng, Bowen Yu, Chang
430	Gao, Chengen Huang, Chenxu Lv, and 1 others.
431	2025. Qwen3 technical report. <i>arXiv preprint</i>
432	<i>arXiv:2505.09388</i> .
433	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali
434	Farhadi, and Yejin Choi. 2019. Hellaswag: Can a
435	machine really finish your sentence? <i>arXiv preprint</i>
436	<i>arXiv:1905.07830</i> .
437	Biao Zhang and Rico Sennrich. 2019. Root mean square
438	layer normalization. <i>Advances in neural information</i>
439	<i>processing systems</i> , 32.

## Appendices

### A Model Details

Our analysis focuses on the Qwen-2.5 model family (Qwen Team, 2025). For base models, we use the same versions selected by DeepSeek-AI (2025): Qwen2.5-Math-1.5B, Qwen2.5-Math-7B, Qwen2.5-14B, and Qwen2.5-32B. For instruction-tuned models, we use the standard instruct variants from the Qwen-2.5 family: Qwen2.5-1.5B-Instruct, Qwen2.5-7B-Instruct, Qwen2.5-14B-Instruct, and Qwen2.5-32B-Instruct. Additionally, we include the corresponding DeepSeek-R1-Distill versions derived from these base models: DeepSeek-R1-Distill-Qwen-1.5B, DeepSeek-R1-Distill-Qwen-7B, DeepSeek-R1-Distill-Qwen-14B, and DeepSeek-R1-Distill-Qwen-32B.

Models of the same size share identical architectures. The architectural details are provided in Table 2.

Table 2: Model details.

Models	Layers	Heads (Q/KV)
1.5B	28	12 / 2
7B	28	28 / 4
14B	48	40 / 8
32B	64	40 / 8

### B Additional Effective Depth Results on GSM8K

We show the results of logit lens in Figure 2, the effects of skipping a layer on future computations in Figure 3(a) and on output distributions in Figure 3(b). Besides, the results of integrated gradients residual erasure are shown in Figure 4(a) and Figure 4(b) respectively.

### C Effective Depth of All Models Evaluated on GSM8K and HellaSwag

We show the results of effective depth of Qwen-2.5 family (base and instruct models) and their long-CoT variants tested on GSM8K and HellaSwag, including residual cosine similarity results in Figure 5; the effects of skipping a layer on future computations in Figure 6 and on output distributions in Figure 7; logit lens KL divergence in Figure 8; logit lens overlap in Figure 9.

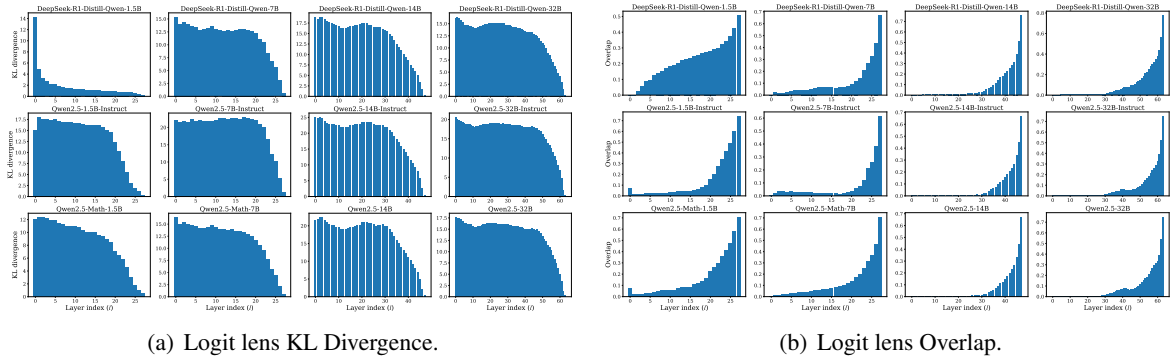
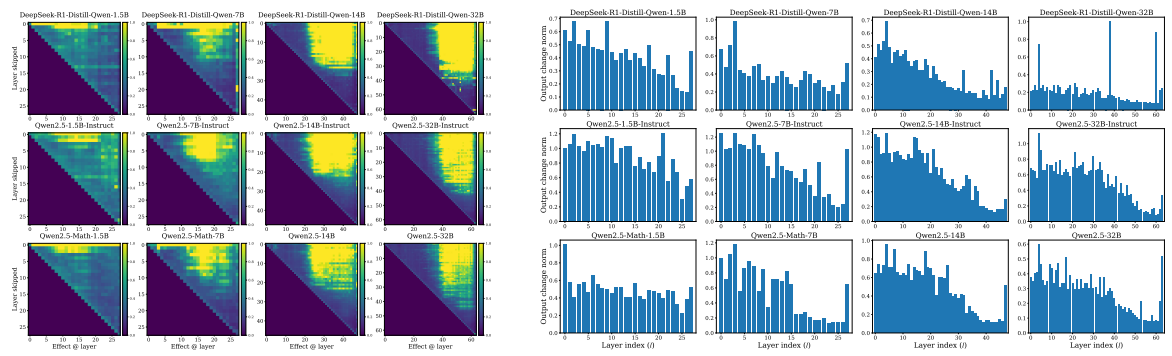
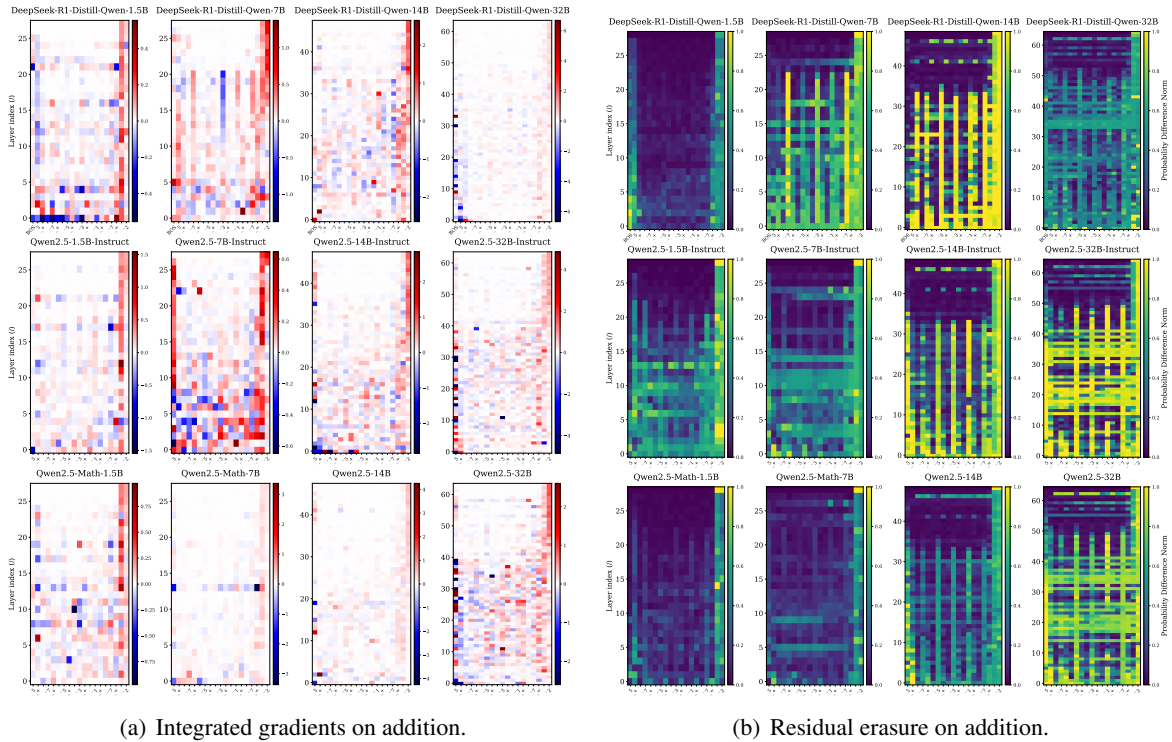


Figure 2: Logit lens Results on GSM8K.



(a) Effects of skipping a layer on later layers' contribution. (b) Effects of skipping a layer on output norm.

Figure 3: Effect of skipping a layer on future computation evaluated on GSM8K.



(a) Integrated gradients on addition. (b) Residual erasure on addition.

Figure 4: The Effects of individual computation steps evaluated on GSM8K.

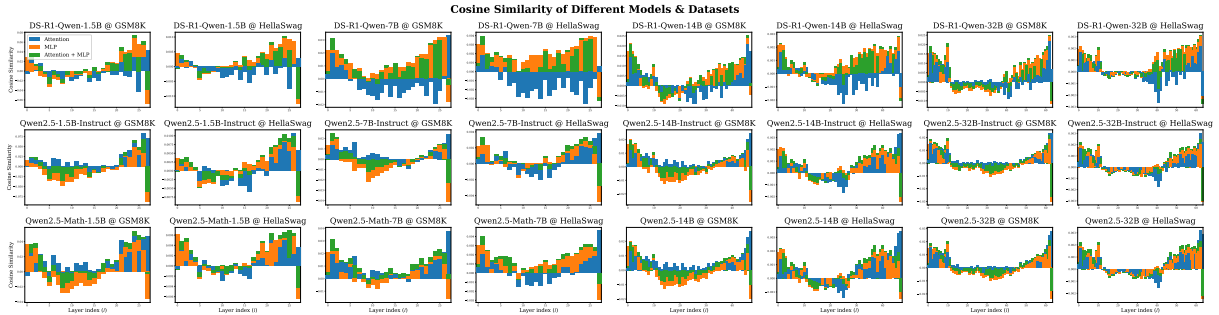


Figure 5: Residual cosine similarity of all models on GSM8K and HellaSwag.

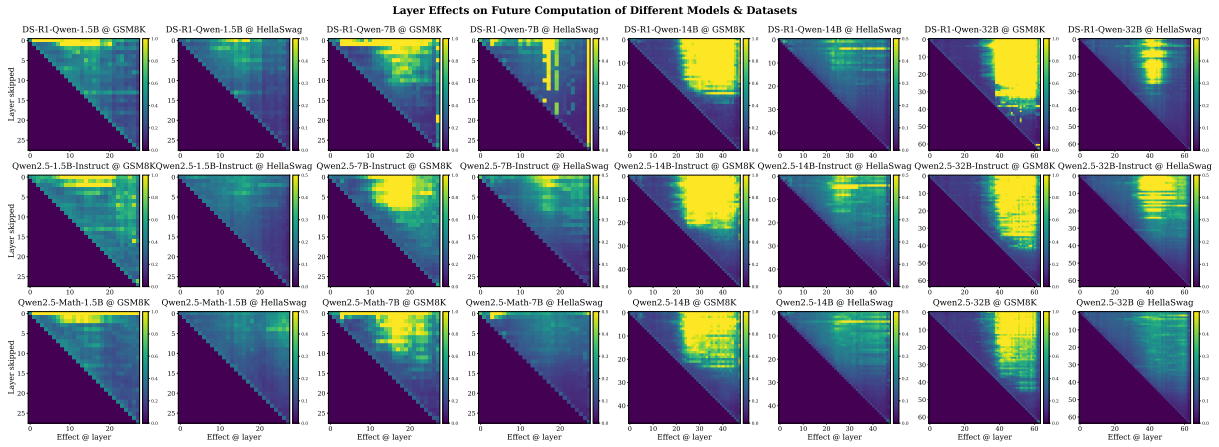


Figure 6: The effects of skipping a layer on future computations, the results include all models on GSM8K and HellaSwag.

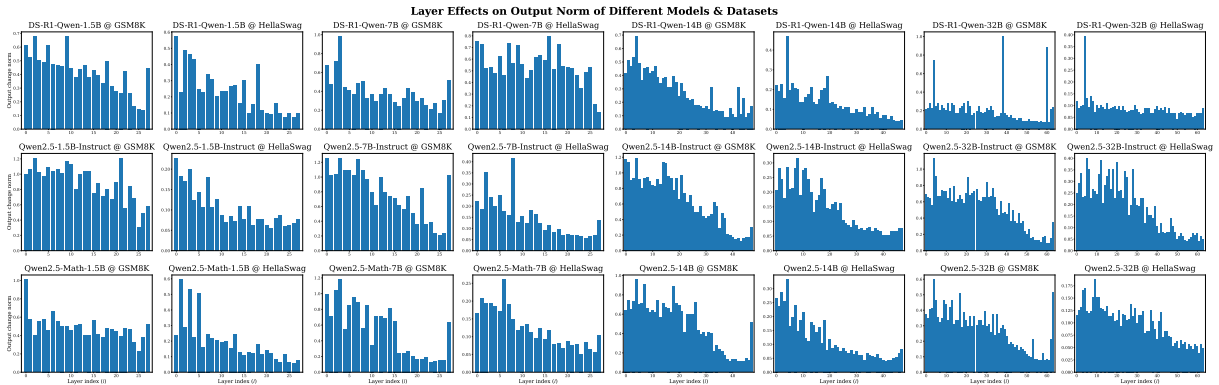


Figure 7: The effects of skipping a layer on output distributions, the results include all models on GSM8K and HellaSwag.

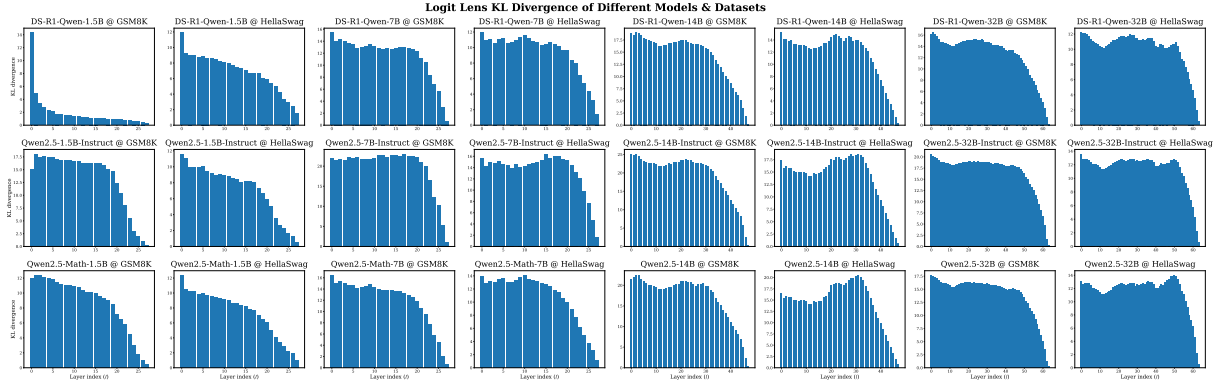


Figure 8: Logit lens KL divergence between early layer distributions and the final distributions. The results include all models on GSM8K and HellaSwag.

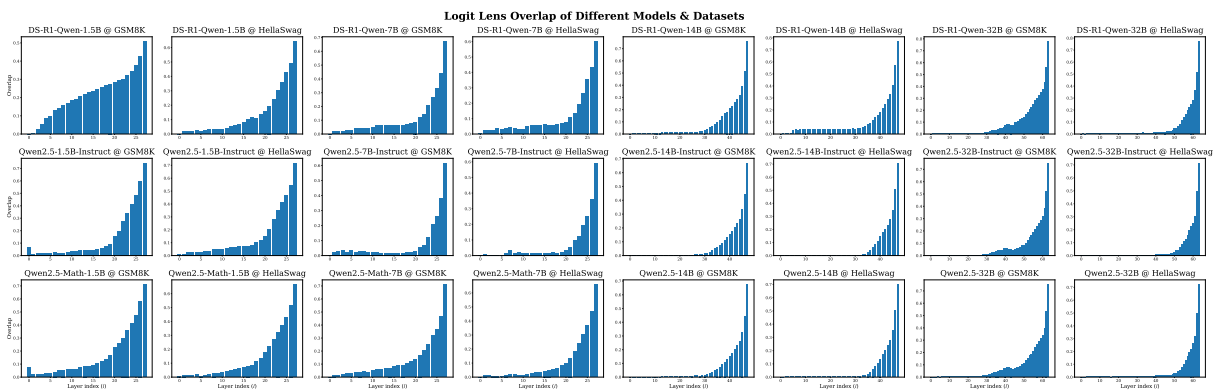


Figure 9: Logit lens top-5 overlap between early layer distributions and the final distributions. The results include all models on GSM8K and HellaSwag.