

Decoupling the Effect of Chain-of-Thought Reasoning: A Human Label Variation Perspective

Anonymous ACL submission

Abstract

Reasoning-tuned LLMs utilizing long Chain-of-Thought (CoT) excel at single-answer tasks, yet their ability to model Human Label Variation—which requires capturing probabilistic ambiguity rather than resolving it—remains underexplored. We investigate this through systematic disentanglement experiments on distribution-based tasks, employing Cross-CoT experiments to isolate the effect of reasoning text from intrinsic model priors. We observe a distinct "decoupled mechanism": while CoT improves distributional alignment, final accuracy is dictated by CoT content (99% variance contribution), whereas distributional ranking is governed by model priors (over 80%). Step-wise analysis further shows that while CoT's influence on accuracy grows monotonically during the reasoning process, distributional structure is largely determined by LLM's intrinsic priors. These findings suggest that long CoT serves as a decisive LLM decision-maker for the top option but fails to function as a granular distribution calibrator for ambiguous tasks.

1 Introduction

Reasoning-tuned large language models (LLMs) with long CoT reasoning achieve strong performance on many benchmarks (Touvron et al., 2023; Dubey et al., 2024; OpenAI, 2023; Wei et al., 2022; Wang et al., 2023; DeepSeek-AI et al., 2025; Team, 2025c; Hurst et al., 2024), usually measured by accuracy under the assumption of a single correct answer (Hendrycks et al., 2021a; Rein et al., 2023; Wang et al., 2024; Sun et al., 2025; Hendrycks et al., 2021b). However, many real-world tasks are inherently ambiguous or subjective, with human annotators often disagreeing due to genuine semantic uncertainty (Pavlick and Kwiatkowski, 2019; Aroyo and Welty, 2015). Such Human Label Variation (HLV) requires models to predict distributions over plausible answers, making argmax-based evaluation insufficient (Uma et al., 2021; Plank, 2022;

Cabitz et al., 2023; Hu et al., 2025). Intuitively, reasoning through intermediate steps might better reflect such variations compared to direct answering (Chen et al., 2025a), motivating us to ask *RQ1: whether long CoT helps models better approximate human label distributions*, and *RQ2: whether any gains come from CoT reasoning or the model's latent parametric knowledge*.

To investigate *RQ1*, we utilize ChaosNLI (Nie et al., 2020), a benchmark capturing collective human opinions. We analyze the latent answer distributions behind CoT using complementary metrics: accuracy for correctness, and Jensen–Shannon Divergence (JSD, Endres and Schindelin 2003) and Spearman's ρ (Spearman, 1961) for distributional and ranking alignment. To further disentangle CoT's role from model-intrinsic priors (*RQ2*), we conduct: i) Cross-CoT experiments, injecting one model's CoT into another to test reasoning transfer; and ii) Step-wise analysis, truncating CoT to track how influence evolves over reasoning steps.

Our analysis uncovers a notable "split influence". While LLMs generally improve distributional alignment (lower JSD) after reasoning, this gain is not uniform across metrics. Using ANOVA to calculate the variance contribution percentage in our Cross-CoT experiments, we find that final accuracy is overwhelmingly determined by the CoT content ($\approx 99\%$), confirming the strong role of reasoning chain in steering the top-1 answer decision. In stark contrast, the distributional structure—ranking and probability allocation among non-argmax options—is largely immune to CoT, remaining governed by model priors ($>80\%$).

Step-wise analysis further clarifies this dynamic. While all metrics evolve throughout reasoning, changes in accuracy are predominantly driven by CoT and grow monotonically with later steps. By comparison, changes in distributional similarity (JSD and Spearman's ρ) are mostly determined by the LLM's intrinsic behavior. This reveals a di-

chotomy: current long CoT paradigms act as strong LLM decision makers but weak distribution calibrators. CoT tends to progressively concentrate probability mass to lock in the most likely answer latently, but fails to govern the reshaping of the probability landscape for alternative options. This work highlights the structural limitations of current reasoning processes in capturing fine-grained answer uncertainty and motivates the need for distribution-aware reasoning mechanisms.

2 Background

HLV in Natural Language Inference. Unlike the single-label assumption in most benchmarks, NLI is often inherently ambiguous: a premise and hypothesis can elicit a spectrum of plausible interpretations, a phenomenon known as HLV (Plank, 2022). Benchmarks such as ChaosNLI capture this by representing labels as probability distributions rather than single gold labels (Nie et al., 2020; Weber-Genzel et al., 2024; Jiang et al., 2023; Hong et al., 2025). Evaluating models under HLV requires moving beyond standard accuracy to distributional metrics that measure alignment with collective human judgments (Kurniawan et al., 2025; Lee et al., 2023; Leonardelli et al., 2023; Chen et al., 2024, 2025b,a; Ni et al., 2025).

Reasoning under Distributional Uncertainty. Recent LLM advancements emphasize reasoning-intensive paradigms. Long CoT enables models to decompose problems into intermediate steps (Wang et al., 2023; DeepSeek-AI et al., 2025; Team, 2025c; Hurst et al., 2024), effectively reducing uncertainty and producing high-confidence conclusions in deterministic tasks. However, its role in probabilistic HLV settings is less clear. Generating explicit reasoning can inadvertently suppress valid alternative interpretations, potentially biasing the model toward the top-1 choice. While prior work has explored confidence-based calibration (Zhao et al., 2025; Yoon et al., 2025; Mao et al., 2025), it remains unclear whether CoT actively shapes the full output distribution or mainly rationalizes the final decision, leaving non-argmax probabilities governed by the model’s intrinsic priors.

3 Experiments

3.1 Setup

Task We experiment on 3 ChaosNLI subsets: MNLI, SNLI, and α NLI (Bowman et al., 2015;

| Reasoning LLMs | Abbr. |
|--|----------|
| Qwen/Qwen3-30B-A3B-Thinking-2507 (Team, 2025b) | Qwen |
| deepseek-ai/DeepSeek-R1-Distill-Llama-70B (DeepSeek-AI et al., 2025) | R1-Llama |
| deepseek-ai/DeepSeek-R1-Distill-Qwen-32B (DeepSeek-AI et al., 2025) | R1-Qwen |
| allenai/Olmo-3-32B-Think (Olmo et al., 2025) | Olmo |
| zai-org/GLM-Z1-32B-0414 (GLM et al., 2024) | GLM |
| ByteDance-Seed/Seed-OSS-36B-Instruct (Team, 2025a) | Seed |
| openai/gpt-oss-20b (OpenAI, 2025) | GPT |

Table 1: Reasoning LLMs and their abbreviation.

Williams et al., 2018; Bhagavatula et al., 2020). Each instance is annotated by 100 crowdworkers, enabling reliable human judgment distributions (HJD). MNLI and SNLI are three-way classification tasks (entailment, neutral, contradiction), yielding 3-d label distributions. α NLI is a binary-choice task, where annotators select the better hypothesis for a given observation pair, producing 2-d distributions.¹ Dataset details are in Appendix A.

Models To comprehensively evaluate the HLV performance of reasoning-tuned LLMs, we select a range of state-of-the-art open-source reasoning models (details in Table 1). All follow a reason-then-answer paradigm: generating a long CoT reasoning process before outputting a final answer.

Evaluation All NLI instances are reformulated as multiple-choice questions. Model predictions are extracted using the first-token probability method (Santurkar et al., 2023; Durmus et al., 2023; Liang et al., 2023), where logits are aggregated and normalized to obtain an output probability distribution over answer options. We measure the HLV alignment between the model-generated distribution and the corresponding HJD using JSD. We also report accuracy. See details in Appendix B.

3.2 Does CoT Improve HLV Performance?

We examine the impact of reasoning by comparing model performance before and after CoT. See Table 2, the effect of CoT on accuracy is mixed. While most models improve on SNLI and α NLI, performance on MNLI is highly unstable. In contrast, JSD consistently decreases across nearly all models and datasets. Importantly, this improved distributional alignment is often decoupled from accuracy: even when the accuracy decreases (e.g., Qwen on MNLI), the output distribution aligns more closely with human judgments on average.

CoT generally reduces JSD, indicating useful signals for HLV, but the benefit varies across models. Models with similar post-CoT accuracy can

¹ChaosNLI is ideal for HLV evaluation as a rare non-social science benchmark with collective HJDs (Hu et al., 2025).

| Task | MNLI | | | | SNLI | | | | α NLI | | | | |
|----------|--------------|---------------------------------|--------------------------------|-----------------------------------|----------------------------------|---------------------------------|--------------------------------|-----------------------------------|----------------------------------|---------------------------------|--------------------------------|-----------------------------------|----------------------------------|
| | LLMs/Metrics | ACC _{start} \uparrow | ACC _{last} \uparrow | JSD _{start} \downarrow | JSD _{last} \downarrow | ACC _{start} \uparrow | ACC _{last} \uparrow | JSD _{start} \downarrow | JSD _{last} \downarrow | ACC _{start} \uparrow | ACC _{last} \uparrow | JSD _{start} \downarrow | JSD _{last} \downarrow |
| Qwen | | 0,688 | 0,644 | 0,093 | 0,080 | 0,668 | 0,778 | 0,144 | 0,119 | 0,749 | 0,890 | 0,108 | 0,084 |
| R1-Llama | | 0,666 | 0,689 | 0,082 | 0,077 | 0,615 | 0,750 | 0,133 | 0,123 | 0,839 | 0,878 | 0,098 | 0,091 |
| R1-Qwen | | 0,734 | 0,672 | 0,080 | 0,072 | 0,689 | 0,764 | 0,127 | 0,115 | 0,832 | 0,860 | 0,094 | 0,081 |
| Olmo | | 0,614 | 0,609 | 0,088 | 0,082 | 0,738 | 0,775 | 0,133 | 0,122 | 0,819 | 0,863 | 0,107 | 0,087 |
| GLM | | 0,670 | 0,640 | 0,082 | 0,077 | 0,545 | 0,756 | 0,134 | 0,120 | 0,834 | 0,888 | 0,099 | 0,088 |
| Seed | | 0,705 | 0,614 | 0,077 | 0,083 | 0,766 | 0,777 | 0,124 | 0,127 | 0,868 | 0,887 | 0,098 | 0,095 |
| GPT | | 0,437 | 0,672 | 0,095 | 0,077 | 0,596 | 0,772 | 0,145 | 0,119 | 0,793 | 0,872 | 0,112 | 0,080 |

Table 2: Results before and after reasoning. *start* and *last* denote before reasoning and after completion. Red indicates an increase, blue a decrease. Arrows next to metric names show whether higher or lower is better.

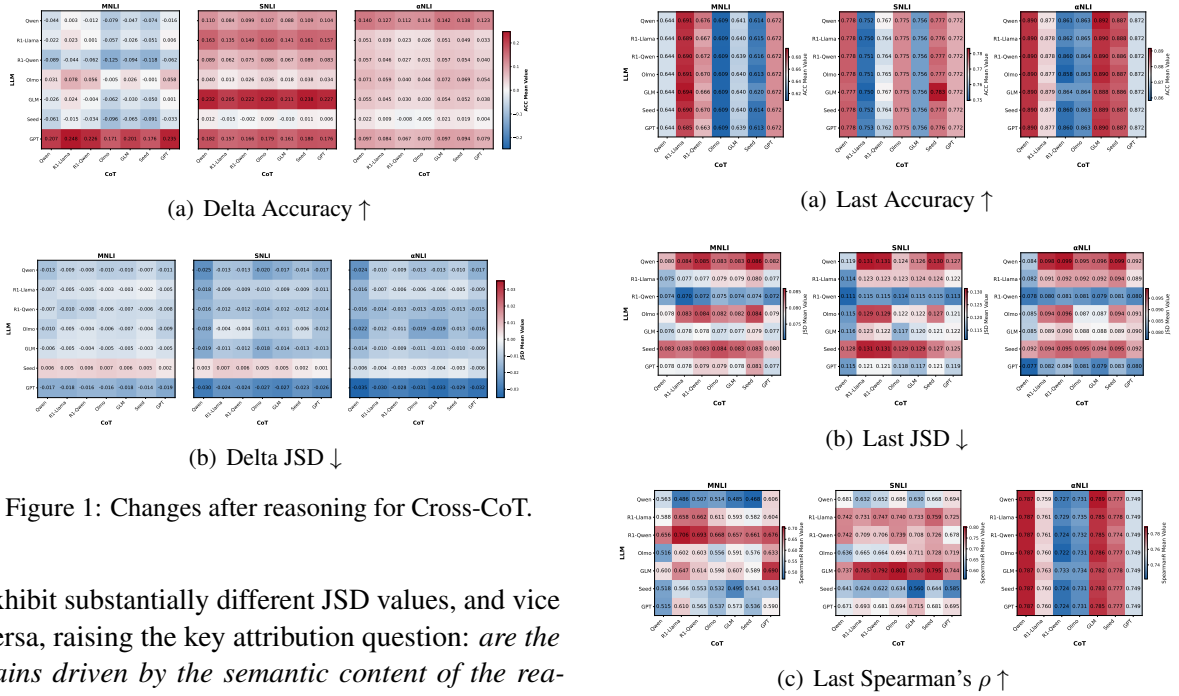


Figure 1: Changes after reasoning for Cross-CoT.

exhibit substantially different JSD values, and vice versa, raising the key attribution question: *are the gains driven by the semantic content of the reasoning itself, or by model-specific inductive biases when interpreting the reasoning text?*

3.3 Cross-CoT Evaluation

To disentangle the source of JSD improvements, we conduct *Cross-CoT* experiments, injecting reasoning paths from different source models into various inference models. We show the performance changes in Figure 1. On MNLI, accuracy shows mixed patterns. In contrast, consistent with single-model results, JSD improves across nearly all Cross-CoT pairings.² Regardless of the reasoning source model, injecting a CoT almost universally reduces divergence from human distributions. This confirms: **CoT text acts as a portable carrier of HLV-relevant information—reasoning generated by one model can facilitate better distributional alignment in another.** However, the divergent patterns between accuracy and JSD motivate us to examine why these two metrics respond differently to CoT, and how CoT influences them.

²The box plot (Figure 4) in Appendix C shows improvements are widespread across instances, not driven by outliers.

Figure 2: Last results after reasoning for Cross-CoT.

4 Analyses

4.1 What Does CoT Determine?

If CoT were the dominant driver, models conditioned on the same CoT should converge to similar outcomes. To test this, we analyze final-step accuracy and JSD. The results (Figure 2(a), 2(b)) reveal a clear dissociation. Accuracy shows a *column-dominant pattern*: for a fixed CoT source, accuracy is nearly identical across inference models, indicating that CoT largely dictates the argmax decision. In contrast, JSD exhibits a *row-dominant pattern*: final divergence is primarily determined by the inference model, with little sensitivity to the CoT.

We further analyze Spearman's ρ as a relaxed non-argmax metric (only rankings). Its heatmap (Figure 2(c)) mirrors the row-dominant structure of JSD rather than accuracy.³ Analysis of Variance (ANOVA) on MNLI confirms this split: CoT

³Note ρ basically equals accuracy on binary-choice α NLI.

| Task | MNL1 | | | | | | | | | SNLI | | | | | | | | | oNLI | | | | | | | | |
|---------|--------|-------|----------|--------|------|----------|-------------------|-------|----------|--------|-------|----------|--------|-------|----------|-------------------|-------|----------|--------|-------|----------|--------|-------|----------|-------------------|-------|----------|
| | ACC | | | JSD | | | Spearman's ρ | | | ACC | | | JSD | | | Spearman's ρ | | | ACC | | | JSD | | | Spearman's ρ | | |
| | LLM | CoT | Residual | LLM | CoT | Residual | LLM | CoT | Residual | LLM | CoT | Residual | LLM | CoT | Residual | LLM | CoT | Residual | LLM | CoT | Residual | LLM | CoT | Residual | LLM | CoT | Residual |
| Step 0 | 100.0% | 0.0% | 0.0% | 100.0% | 0.0% | 0.0% | 100.0% | 0.0% | 0.0% | 100.0% | 0.0% | 0.0% | 100.0% | 0.0% | 0.0% | 100.0% | 0.0% | 0.0% | 100.0% | 0.0% | 0.0% | 100.0% | 0.0% | 0.0% | 100.0% | 0.0% | 0.0% |
| Step 1 | 97.3% | 0.2% | 2.5% | 96.9% | 0.4% | 2.7% | 92.7% | 0.4% | 6.9% | 96.5% | 1.4% | 2.2% | 94.4% | 0.9% | 4.7% | 83.2% | 2.4% | 14.4% | 93.9% | 1.4% | 4.7% | 95.4% | 1.3% | 3.3% | 94.0% | 1.3% | 4.8% |
| Step 2 | 95.7% | 0.5% | 3.7% | 95.3% | 0.7% | 4.0% | 92.3% | 1.8% | 5.9% | 85.1% | 8.6% | 6.4% | 88.2% | 4.4% | 7.4% | 76.1% | 8.0% | 15.8% | 75.7% | 10.5% | 13.8% | 89.3% | 3.9% | 6.8% | 76.3% | 10.4% | 13.2% |
| Step 3 | 91.8% | 1.6% | 6.6% | 93.4% | 1.6% | 5.0% | 91.6% | 2.6% | 5.8% | 75.1% | 16.7% | 8.1% | 82.2% | 8.3% | 9.5% | 73.2% | 12.2% | 14.6% | 64.1% | 22.5% | 13.7% | 81.4% | 8.8% | 9.8% | 64.3% | 22.4% | 13.3% |
| Step 4 | 89.1% | 2.2% | 8.7% | 91.5% | 2.9% | 5.6% | 90.5% | 3.6% | 5.9% | 62.5% | 29.1% | 8.5% | 75.4% | 14.0% | 10.6% | 69.9% | 16.4% | 13.7% | 52.1% | 31.2% | 16.8% | 80.8% | 9.5% | 9.6% | 51.9% | 32.6% | 15.5% |
| Step 5 | 85.2% | 4.3% | 10.6% | 89.3% | 4.2% | 6.4% | 88.0% | 4.9% | 7.1% | 57.4% | 32.4% | 10.2% | 72.8% | 16.6% | 10.6% | 70.8% | 13.0% | 16.2% | 45.2% | 41.7% | 13.1% | 81.6% | 10.4% | 8.0% | 47.2% | 41.3% | 11.5% |
| Step 6 | 79.4% | 8.5% | 12.0% | 87.9% | 5.7% | 6.5% | 83.6% | 7.6% | 8.8% | 58.9% | 29.7% | 11.3% | 72.5% | 17.4% | 10.1% | 71.9% | 10.4% | 17.7% | 35.4% | 53.8% | 10.8% | 84.9% | 9.4% | 5.8% | 36.4% | 53.3% | 10.3% |
| Step 7 | 66.5% | 18.4% | 15.1% | 87.5% | 7.1% | 6.4% | 77.5% | 10.1% | 12.4% | 53.4% | 30.6% | 16.0% | 73.1% | 17.2% | 9.7% | 73.9% | 6.6% | 19.5% | 22.5% | 67.3% | 10.2% | 87.0% | 8.2% | 4.9% | 22.1% | 67.5% | 10.4% |
| Step 8 | 44.1% | 38.3% | 17.5% | 86.4% | 7.0% | 6.7% | 73.3% | 14.5% | 12.2% | 43.0% | 33.2% | 23.8% | 72.1% | 19.1% | 8.8% | 71.8% | 7.9% | 20.3% | 9.6% | 76.3% | 14.2% | 85.2% | 10.1% | 4.7% | 10.1% | 76.9% | 13.0% |
| Step 9 | 22.4% | 65.2% | 12.4% | 84.3% | 8.5% | 7.2% | 73.0% | 16.8% | 10.3% | 34.8% | 43.9% | 21.2% | 71.3% | 19.8% | 8.9% | 77.8% | 6.3% | 15.9% | 5.3% | 83.5% | 11.2% | 81.4% | 13.2% | 5.4% | 5.4% | 84.1% | 10.5% |
| Step 10 | 0.1% | 99.5% | 0.4% | 83.3% | 8.7% | 8.1% | 71.1% | 13.1% | 15.7% | 0.2% | 98.6% | 1.2% | 67.7% | 22.2% | 10.1% | 81.7% | 3.3% | 15.0% | 0.1% | 99.4% | 0.5% | 76.9% | 15.2% | 7.9% | 0.1% | 99.4% | 0.5% |

Table 3: Step-wise ANOVA results. Each CoT is split into 10 segments by sentence, yielding 11 intermediate answers from no-thinking (step 0) to full-thinking (step 10). ANOVA is computed for each Cross-CoT heatmap. Red numbers indicate the factor dominating the metric at that step. All step-wise heatmaps are in Appendix F.

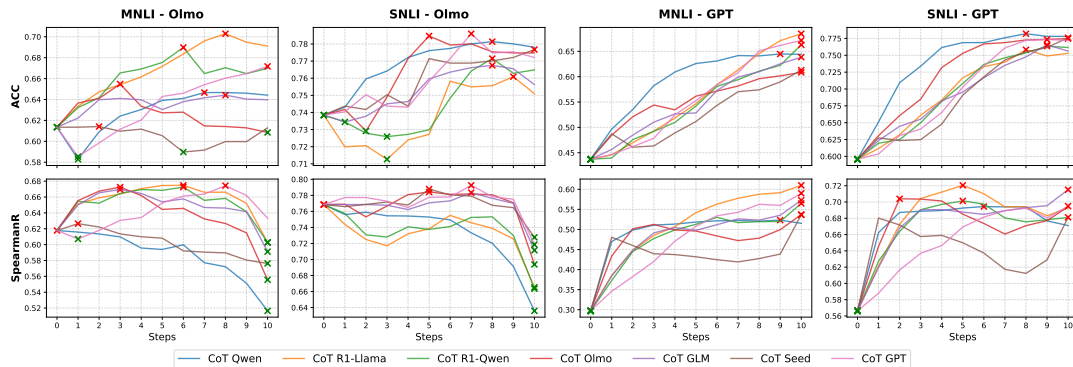


Figure 3: Curve cases for step-wise evaluation. Max and min points are marked. All results are in Appendix G.

explains 99% of the variance in accuracy, but only a small fraction in JSD and Spearman's ρ (8.7% and 13.1%), where model identity dominates (83.3% and 71.1%).⁴ This asymmetry exposes a fundamental limitation of current CoT paradigms. CoT is highly effective at explicit decision-making, capable of overriding a model's prior to determine the argmax. However, in the non-argmax space—namely, the ranking and probability allocation over alternative options—its influence sharply diminishes. **Models appear to follow CoT for the final choice, but revert to their latent parametric preferences when distributing uncertainty.**

4.2 When Does CoT Take Control?

Although CoT explains nearly all variance in accuracy, its impact on distributional metrics remains limited. To understand how this asymmetry develops, we apply early stopping to the CoT, truncating it at fixed increments and evaluating intermediate performance.⁵ Step-wise ANOVA (Table 3) reveals a sharp divergence. CoT influence on accuracy remains modest during reasoning, then spikes abruptly at the final step, forming a clear inflection. In contrast, its influence on JSD and Spearman's ρ stays uniformly low, with no point at which CoT

⁴Details of ANOVA are in Appendix D.

⁵Implementation Details are in the Appendix E. We will release all codes, CoTs, logits upon publication.

overrides the model's ranking behavior.

This pattern is also illustrated by representative models (Figure 3). Accuracy often shifts or converges only at the conclusion, while Spearman's ρ fluctuates without a consistent trend. Thus, CoT determines the LLM's final choice but not the structure of uncertainty. Our anecdotal evidence in Appendix H supports that this can be attributed to the CoT format: standard CoT often ends with an explicit conclusion, providing a strong argmax signal, while distributional cues remain implicit. Consequently, **models leverage CoT for decision-making but revert to intrinsic priors for probability allocation, exposing a structural inability of raw CoT to shape answer distributions.**

5 Conclusion

From an HLV perspective, we identify a fundamental "split influence" in CoT reasoning: while reasoning content predominantly determines the LLM's final argmax choice latently, the probability landscape of alternative options remains anchored to the model's intrinsic priors. This exposes a structural limitation where standard CoT effectively collapses ambiguity for decision-making but fails to calibrate fine-grained uncertainty for alternative, plausible answers. Consequently, advancing HLV modeling requires moving beyond implicit reasoning traces toward distribution-aware paradigms.

266 Limitations

267 Our work has two main limitations. First, our eval-
268 uation relies solely on final human label distribu-
269 tions, as ChaosNLI lacks annotated intermediate
270 reasoning steps. Consequently, our step-wise anal-
271 ysis compares intermediate model outputs against
272 the final human consensus rather than step-specific
273 ground truth. Addressing this limitation would re-
274 quire future improvements in human annotation,
275 where reasoning steps and intermediate answers
276 are collected for direct comparison. An alternative
277 approach could be the use of relative references:
278 for example, treating the intermediate answers gen-
279 erated by a CoT-provider LLM as the gold standard
280 to evaluate the faithfulness of other LLMs. An-
281 other possibility is to employ an entailment model
282 to determine whether each reasoning step in a CoT
283 entails the previous step, thereby inferring interme-
284 diate answers recursively. However, both of these
285 approaches are highly dependent on the accuracy
286 of the model itself; inaccuracies could introduce
287 evaluation biases.

288 Second, our study does not include a direct hu-
289 man evaluation of the textual content of the CoTs.
290 Instead, we focus on assessing CoTs in terms of
291 their impact on LLM behavior, especially answer
292 distributions. While this approach emphasizes the
293 effect of reasoning on model outputs, it overlooks
294 the quality of the specific text content. Conduct-
295 ing human evaluation of long CoTs is particularly
296 resource-intensive, given their length and the effort
297 required to annotate individual sentences and their
298 interrelations. Nevertheless, considering the grow-
299 ing importance of CoT reasoning in NLP, carefully
300 designed human evaluation to verify whether ex-
301 treme values in reasoning metrics correspond to rea-
302 sonable positions in the text represents a promising
303 direction for future work. Such evaluation could
304 help us better understand the effects of CoT on
305 LLM reasoning.

306 Ethical Considerations

307 This work primarily involves the analysis of NLI
308 datasets and open-sourced LLMs. All data used
309 are publicly available and do not contain person-
310 ally identifiable information. No sensitive or po-
311 tentially harmful content is generated or utilized
312 in this study. Therefore, we do not anticipate any
313 ethical concerns arising from our work.

Use of AI Assistants The authors acknowledge
the use of ChatGPT solely for correcting grammat-
ical errors, enhancing the coherence of the final
manuscript.

References

- Lora Aroyo and Chris Welty. 2015. [Truth is a lie: Crowd truth and the seven myths of human annotation](#). *AI Mag.*, 36(1):15–24.
- Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Han-nah Rashkin, Doug Downey, Wen-tau Yih, and Yejin Choi. 2020. [Abductive commonsense reasoning](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Federico Cabitza, Andrea Campagner, and Valerio Basile. 2023. [Toward a perspectivist turn in ground truthing for predictive computing](#). In *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 6860–6868. AAAI Press.
- Beiduo Chen, Yang Janet Liu, Anna Korhonen, and Barbara Plank. 2025a. [Threading the needle: Reweaving chain-of-thought reasoning to explain human label variation](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 33099–33123, Suzhou, China. Association for Computational Linguistics.
- Beiduo Chen, Siyao Peng, Anna Korhonen, and Barbara Plank. 2025b. [A rose by any other name: LLM-generated explanations are good proxies for human explanations to collect label distributions on NLI](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 10777–10802, Vienna, Austria. Association for Computational Linguistics.
- Beiduo Chen, Xinpeng Wang, Siyao Peng, Robert Litschko, Anna Korhonen, and Barbara Plank. 2024. [“seeing the big through the small”: Can LLMs approximate human judgment distributions on NLI from a few explanations?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 14396–14419, Miami, Florida, USA. Association for Computational Linguistics.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,

| | | | |
|-----|---|--|-----|
| 369 | Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, | Simbench: Benchmarking the ability of large language models to simulate human behaviors. | 426 |
| 370 | Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi- | <i>CoRR</i> , abs/2510.17516. | 427 |
| 371 | hong Shao, Zhuoshu Li, Ziyi Gao, and 81 others. | | 428 |
| 372 | 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. | | |
| 373 | <i>CoRR</i> , abs/2501.12948. | | |
| 374 | | | |
| 375 | Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, | Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam | 429 |
| 376 | Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, | Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, | 430 |
| 377 | Akhil Mathur, Alan Schelten, Amy Yang, Angela | Akila Welihinda, Alan Hayes, Alec Radford, Alek- | 431 |
| 378 | Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, | sander Madry, Alex Baker-Whitcomb, Alex Beutel, | 432 |
| 379 | Archi Mitra, Archie Sravankumar, Artem Korenev, | Alex Borzunov, Alex Carney, Alex Chow, Alex Kir- | 433 |
| 380 | Arthur Hinsvark, Arun Rao, Aston Zhang, and 82 | illov, Alex Nichol, Alex Paino, and 79 others. 2024. | 434 |
| 381 | others. 2024. The Llama 3 Herd of Models. | GPT-4o System Card. | 435 |
| 382 | <i>CoRR</i> , abs/2407.21783. | | |
| 383 | Esin Durmus, Karina Nyugen, Thomas I. Liao, Nicholas | Nan-Jiang Jiang, Chenhao Tan, and Marie-Catherine | 436 |
| 384 | Schiefer, Amanda Askell, Anton Bakhtin and Abhi- | de Marneffe. 2023. Ecologically valid explanations | 437 |
| 385 | manyu Carol Chen, Zac Hatfield-Dodds, Danny Her- | for label variation in NLI. In <i>Findings of the As-</i> | 438 |
| 386 | nanandez, Nicholas Joseph, Liane Lovitt, Sam McCan- | <i>sociation for Computational Linguistics: EMNLP</i> | 439 |
| 387 | dlish, Orowa Sikder, Alex Tamkin, Janel Thamkul, | 2023, pages 10622–10633, Singapore. Association | 440 |
| 388 | Jared Kaplan, Jack Clark, and Deep Ganguli. 2023. | for Computational Linguistics. | 441 |
| 389 | Towards measuring the representation of subjective | | |
| 390 | global opinions in language models. | Kemal Kurniawan, Meladel Mistica, Timothy Baldwin, | 442 |
| 391 | <i>CoRR</i> , abs/2306.16388. | and Jey Han Lau. 2025. Training and evaluating with | 443 |
| 392 | | human label variation: An empirical study. | 444 |
| 393 | Dominik Maria Endres and Johannes E. Schindelin. | <i>CoRR</i> , abs/2502.01891. | 445 |
| 394 | 2003. A new metric for probability distributions. | | |
| 395 | <i>IEEE Trans. Inf. Theory</i> , 49(7):1858–1860. | Noah Lee, Na Min An, and James Thorne. 2023. Can | 446 |
| 396 | Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chen- | large language models capture dissenting human | 447 |
| 397 | hui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Han- | voices? In <i>Proceedings of the 2023 Conference</i> | 448 |
| 398 | lin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai | <i>on Empirical Methods in Natural Language Process-</i> | 449 |
| 399 | Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, | <i>ing</i> , pages 4569–4585, Singapore. Association for | 450 |
| 400 | Jing Zhang, Juanzi Li, and 37 others. 2024. Chatglm: | Computational Linguistics. | 451 |
| 401 | A family of large language models from glm-130b to | | |
| 402 | glm-4 all tools. | Elisa Leonardelli, Gavin Abercrombie, Dina Almanea, | 452 |
| 403 | <i>Preprint</i> , arXiv:2406.12793. | Valerio Basile, Tommaso Fornaciari, Barbara Plank, | 453 |
| 404 | | Verena Rieser, Alexandra Uma, and Massimo Poe- | 454 |
| 405 | Dan Hendrycks, Collin Burns, Steven Basart, Andy | sio. 2023. SemEval-2023 task 11: Learning with | 455 |
| 406 | Zou, Mantas Mazeika, Dawn Song, and Jacob Stein- | disagreements (LeWiDi). In <i>Proceedings of the</i> | 456 |
| 407 | hardt. 2021a. Measuring massive multitask language | <i>17th International Workshop on Semantic Evaluation</i> | 457 |
| 408 | understanding. In <i>9th International Conference on</i> | <i>(SemEval-2023)</i> , pages 2304–2318, Toronto, Canada. | 458 |
| 409 | <i>Learning Representations, ICLR 2021, Virtual Event,</i> | Association for Computational Linguistics. | 459 |
| 410 | <i>Austria, May 3-7, 2021.</i> | | |
| 411 | OpenReview.net. | Percy Liang, Rishi Bommasani, Tony Lee, Dimitris | 460 |
| 412 | Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul | Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian | 461 |
| 413 | Arora, Steven Basart, Eric Tang, Dawn Song, and | Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Ku- | 462 |
| 414 | Jacob Steinhardt. 2021b. Measuring mathematical | mar, Benjamin Newman, Binhang Yuan, Bobby Yan, | 463 |
| 415 | problem solving with the MATH dataset. In <i>Pro-</i> | Ce Zhang, Christian Cosgrove, Christopher D. Man- | 464 |
| 416 | <i>ceedings of the Neural Information Processing Sys-</i> | ning, Christopher Ré, Diana Acosta-Navas, Drew A. | 465 |
| 417 | <i>tems Track on Datasets and Benchmarks 1, NeurIPS</i> | Hudson, and 31 others. 2023. Holistic evaluation of | 466 |
| 418 | <i>Datasets and Benchmarks 2021, December 2021, vir-</i> | language models. | 467 |
| 419 | <i>tual.</i> | <i>Trans. Mach. Learn. Res.</i> , 2023. | 468 |
| 420 | | Zhenjiang Mao, Artem Bisliouk, Rohith Reddy Nama, | 469 |
| 421 | Pingjun Hong, Beiduo Chen, Siyao Peng, Marie- | and Ivan Ruchkin. 2025. Temporalizing confidence: | 470 |
| 422 | Catherine de Marneffe, and Barbara Plank. 2025. | Evaluation of chain-of-thought reasoning with signal | 471 |
| 423 | LiTeX: A linguistic taxonomy of explanations for un- | temporal logic. | 472 |
| 424 | derstanding within-label variation in natural language | | |
| 425 | inference. In <i>Proceedings of the 2025 Conference on</i> | Jingwei Ni, Yu Fan, Vilém Zouhar, Donya Rooein, | 473 |
| 426 | <i>Empirical Methods in Natural Language Processing,</i> | Alexander Hoyle, Mrinmaya Sachan, Markus Leip- | 474 |
| 427 | pages 34053–34073, Suzhou, China. Association for | pold, Dirk Hovy, and Elliott Ash. 2025. Can reason- | 475 |
| 428 | Computational Linguistics. | ing help large language models capture human anno- | 476 |
| 429 | | tator disagreement? | 477 |
| 430 | Tiancheng Hu, Joachim Baumann, Lorenzo Lupo, | <i>Preprint</i> , arXiv:2506.19467. | 478 |
| 431 | Nigel Collier, Dirk Hovy, and Paul Röttger. 2025. | | 479 |
| 432 | | Yixin Nie, Xiang Zhou, and Mohit Bansal. 2020. What | 480 |
| 433 | | can we learn from collective human opinions on nat- | 481 |
| 434 | | ural language inference data? In <i>Proceedings of the</i> | 482 |
| 435 | | <i>2020 Conference on Empirical Methods in Natural</i> | |
| 436 | | <i>Language Processing (EMNLP)</i> , pages 9131–9143, | |
| 437 | | Online. Association for Computational Linguistics. | |

| | | |
|-----|---|-----|
| 483 | Team Olmo, Allyson Ettinger, Amanda Bertsch, Bailey Kuehl, David Graham, David Heineman, Dirk Groeneveld, Faeze Brahman, Finbarr Timbers, Hamish Ivison, and 1 others. 2025. Olmo 3. <i>arXiv preprint arXiv:2512.13961</i> . | 537 |
| 484 | | 538 |
| 485 | | 539 |
| 486 | | 540 |
| 487 | | |
| 488 | OpenAI. 2023. GPT-4 technical report . <i>CoRR</i> , abs/2303.08774. | 541 |
| 489 | | 542 |
| 490 | OpenAI. 2025. gpt-oss-120b & gpt-oss-20b model card . <i>Preprint</i> , arXiv:2508.10925. | 543 |
| 491 | | 544 |
| 492 | Ellie Pavlick and Tom Kwiatkowski. 2019. Inherent disagreements in human textual inferences . <i>Transactions of the Association for Computational Linguistics</i> , 7:677–694. | 545 |
| 493 | | 546 |
| 494 | | 547 |
| 495 | | |
| 496 | Barbara Plank. 2022. The “problem” of human label variation: On ground truth in data, modeling and evaluation . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 10671–10682, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. | 548 |
| 497 | | 549 |
| 498 | | 550 |
| 499 | | 551 |
| 500 | | 552 |
| 501 | | 553 |
| 502 | David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2023. GPQA: A graduate-level google-proof q&a benchmark . <i>CoRR</i> , abs/2311.12022. | 554 |
| 503 | | 555 |
| 504 | | 556 |
| 505 | | 557 |
| 506 | | |
| 507 | Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose opinions do language models reflect? In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 29971–30004. PMLR. | 558 |
| 508 | | 559 |
| 509 | | 560 |
| 510 | | 561 |
| 511 | | 562 |
| 512 | | 563 |
| 513 | | 564 |
| 514 | Charles Spearman. 1961. The proof and measurement of association between two things. | 565 |
| 515 | | 566 |
| 516 | Haoxiang Sun, Yingqian Min, Zhipeng Chen, Wayne Xin Zhao, Zheng Liu, Zhongyuan Wang, Lei Fang, and Ji-Rong Wen. 2025. Challenging the boundaries of reasoning: An olympiad-level math benchmark for large language models . <i>CoRR</i> , abs/2503.21380. | 567 |
| 517 | | 568 |
| 518 | | 569 |
| 519 | | 570 |
| 520 | | 571 |
| 521 | | 572 |
| 522 | ByteDance Seed Team. 2025a. Seed-oss open-source models. https://github.com/ByteDance-Seed/seed-oss . | 573 |
| 523 | | 574 |
| 524 | | 575 |
| 525 | Qwen Team. 2025b. Qwen3 technical report . <i>Preprint</i> , arXiv:2505.09388. | 576 |
| 526 | | 577 |
| 527 | | 578 |
| 528 | | 579 |
| 529 | Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>CoRR</i> , abs/2307.09288. | 580 |
| 530 | | 581 |
| 531 | | 582 |
| 532 | | 583 |
| 533 | | 584 |
| 534 | | 585 |
| 535 | | 586 |
| 536 | | |
| | Alexandra Uma, Tommaso Fornaciari, Dirk Hovy, Silviu Paun, Barbara Plank, and Massimo Poesio. 2021. Learning from disagreement: A survey . <i>J. Artif. Intell. Res.</i> , 72:1385–1470. | 587 |
| | | 588 |
| | | 589 |
| | | 590 |
| | | 591 |
| | Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net. | 592 |
| | | 593 |
| | | 594 |
| | | 595 |
| | | 596 |
| | | 597 |
| | | 598 |
| | | 599 |
| | | 600 |
| | | 601 |
| | | 602 |
| | | 603 |
| | | 604 |
| | | 605 |
| | | 606 |
| | | 607 |
| | | 608 |
| | | 609 |
| | | 610 |
| | | 611 |
| | | 612 |
| | | 613 |
| | | 614 |
| | | 615 |
| | | 616 |
| | | 617 |
| | | 618 |
| | | 619 |
| | | 620 |
| | | 621 |
| | | 622 |
| | | 623 |
| | | 624 |
| | | 625 |
| | | 626 |
| | | 627 |
| | | 628 |
| | | 629 |
| | | 630 |
| | | 631 |
| | | 632 |
| | | 633 |
| | | 634 |
| | | 635 |
| | | 636 |
| | | 637 |
| | | 638 |
| | | 639 |
| | | 640 |
| | | 641 |
| | | 642 |
| | | 643 |
| | | 644 |
| | | 645 |
| | | 646 |
| | | 647 |
| | | 648 |
| | | 649 |
| | | 650 |
| | | 651 |
| | | 652 |
| | | 653 |
| | | 654 |
| | | 655 |
| | | 656 |
| | | 657 |
| | | 658 |
| | | 659 |
| | | 660 |
| | | 661 |
| | | 662 |
| | | 663 |
| | | 664 |
| | | 665 |
| | | 666 |
| | | 667 |
| | | 668 |
| | | 669 |
| | | 670 |
| | | 671 |
| | | 672 |
| | | 673 |
| | | 674 |
| | | 675 |
| | | 676 |
| | | 677 |
| | | 678 |
| | | 679 |
| | | 680 |
| | | 681 |
| | | 682 |
| | | 683 |
| | | 684 |
| | | 685 |
| | | 686 |
| | | 687 |
| | | 688 |
| | | 689 |
| | | 690 |
| | | 691 |
| | | 692 |
| | | 693 |
| | | 694 |
| | | 695 |
| | | 696 |
| | | 697 |
| | | 698 |
| | | 699 |
| | | 700 |
| | | 701 |
| | | 702 |
| | | 703 |
| | | 704 |
| | | 705 |
| | | 706 |
| | | 707 |
| | | 708 |
| | | 709 |
| | | 710 |
| | | 711 |
| | | 712 |
| | | 713 |
| | | 714 |
| | | 715 |
| | | 716 |
| | | 717 |
| | | 718 |
| | | 719 |
| | | 720 |
| | | 721 |
| | | 722 |
| | | 723 |
| | | 724 |
| | | 725 |
| | | 726 |
| | | 727 |
| | | 728 |
| | | 729 |
| | | 730 |
| | | 731 |
| | | 732 |
| | | 733 |
| | | 734 |
| | | 735 |
| | | 736 |
| | | 737 |
| | | 738 |
| | | 739 |
| | | 740 |
| | | 741 |
| | | 742 |
| | | 743 |
| | | 744 |
| | | 745 |
| | | 746 |
| | | 747 |
| | | 748 |
| | | 749 |
| | | 750 |
| | | 751 |
| | | 752 |
| | | 753 |
| | | 754 |
| | | 755 |
| | | 756 |
| | | 757 |
| | | 758 |
| | | 759 |
| | | 760 |
| | | 761 |
| | | 762 |
| | | 763 |
| | | 764 |
| | | 765 |
| | | 766 |
| | | 767 |
| | | 768 |
| | | 769 |
| | | 770 |
| | | 771 |
| | | 772 |
| | | 773 |
| | | 774 |
| | | 775 |
| | | 776 |
| | | 777 |
| | | 778 |
| | | 779 |
| | | 780 |
| | | 781 |
| | | 782 |
| | | 783 |
| | | 784 |
| | | 785 |
| | | 786 |
| | | 787 |
| | | 788 |
| | | 789 |
| | | 790 |
| | | 791 |
| | | 792 |
| | | 793 |
| | | 794 |
| | | 795 |
| | | 796 |
| | | 797 |
| | | 798 |
| | | 799 |
| | | 800 |
| | | 801 |
| | | 802 |
| | | 803 |
| | | 804 |
| | | 805 |
| | | 806 |
| | | 807 |
| | | 808 |
| | | 809 |
| | | 810 |
| | | 811 |
| | | 812 |
| | | 813 |
| | | 814 |
| | | 815 |
| | | 816 |
| | | 817 |
| | | 818 |
| | | 819 |
| | | 820 |
| | | 821 |
| | | 822 |
| | | 823 |
| | | 824 |
| | | 825 |
| | | 826 |
| | | 827 |
| | | 828 |
| | | 829 |
| | | 830 |
| | | 831 |
| | | 832 |
| | | 833 |
| | | 834 |
| | | 835 |
| | | 836 |
| | | 837 |
| | | 838 |
| | | 839 |
| | | 840 |
| | | 841 |
| | | 842 |
| | | 843 |
| | | 844 |
| | | 845 |
| | | 846 |
| | | 847 |
| | | 848 |
| | | 849 |
| | | 850 |
| | | 851 |
| | | 852 |
| | | 853 |
| | | 854 |
| | | 855 |
| | | 856 |
| | | 857 |
| | | 858 |
| | | 859 |
| | | 860 |
| | | 861 |
| | | 862 |
| | | 863 |
| | | 864 |
| | | 865 |
| | | 866 |
| | | 867 |
| | | 868 |
| | | 869 |
| | | 870 |
| | | 871 |
| | | 872 |
| | | 873 |
| | | 874 |
| | | 875 |
| | | 876 |
| | | 877 |
| | | 878 |
| | | 879 |
| | | 880 |
| | | 881 |
| | | 882 |
| | | 883 |
| | | 884 |
| | | 885 |
| | | 886 |
| | | 887 |
| | | 888 |
| | | 889 |
| | | 890 |
| | | 891 |
| | | 892 |
| | | 893 |
| | | 894 |
| | | 895 |
| | | 896 |
| | | 897 |
| | | 898 |
| | | 899 |
| | | 900 |
| | | 901 |
| | | 902 |
| | | 903 |
| | | 904 |
| | | 905 |
| | | 906 |
| | | 907 |
| | | 908 |
| | | 909 |
| | | 910 |
| | | 911 |
| | | 912 |
| | | 913 |
| | | 914 |
| | | 915 |
| | | 916 |
| | | 917 |
| | | 918 |
| | | 919 |
| | | 920 |
| | | 921 |
| | | 922 |
| | | 923 |
| | | 924 |
| | | 925 |
| | | 926 |
| | | 927 |
| | | 928 |
| | | 929 |
| | | 930 |
| | | 931 |
| | | 932 |
| | | 933 |
| | | 934 |
| | | 935 |
| | | 936 |
| | | 937 |
| | | 938 |
| | | 939 |
| | | 940 |
| | | 941 |
| | | 942 |
| | | 943 |
| | | 944 |
| | | 945 |
| | | 946 |
| | | 947 |
| | | 948 |
| | | 949 |
| | | 950 |
| | | 951 |
| | | 952 |
| | | 953 |
| | | 954 |
| | | 955 |
| | | 956 |
| | | 957 |
| </ | | |

A Datasets

To evaluate the model’s ability to capture collective human uncertainty and label disagreement, we utilize the **ChaosNLI** dataset (Nie et al., 2020). Unlike standard NLI benchmarks that typically rely on a single “gold” label derived from a majority vote (often among 3–5 annotators), ChaosNLI provides a dense distribution of human annotations.

- **Data Source:** The dataset consists of purely English examples, selected from of SNLI (1514 items, Bowman et al. 2015), MNLI (1599 items, Williams et al. 2018) and α NLI (1532 items, Bhagavatula et al. 2020).
- **Selection Criteria:** The examples were specifically chosen to target ambiguous instances. The authors filtered for examples where the original annotators disagreed (e.g., a 3 vs. 2 vote split) or where the model predictions significantly deviated from the majority label.
- **Annotation Process:** Each example in ChaosNLI is annotated by a crowd of $N = 100$ independent workers. This high volume of annotators allows for the estimation of a true label distribution y_{human} over the three classes (Entailment, Neutral, Contradiction), rather than a deterministic class label. As for α NLI, annotators are asked to select the better hypothesis from a sentence pair for a given observation pair.
- **Objective:** The dataset serves as a testbed for measuring how well a model’s predicted probability distribution p_{model} aligns with the distribution of human judgment y_{human} , often measured via Jensen-Shannon Divergence (JSD).

B Evaluation Details

This section elaborates on the details of the experimental setup. We first describe the experiment details, and then introduce how the NLI task is transformed into a multiple-choice question answering (MCQA) format. We finally introduce the procedure for extracting and converting first-token probabilities, followed by a formal definition of the evaluation metrics used in this paper.

B.1 Experiment Details

All LLMs are evaluated using the initial or recommended parameter settings provided by their respective developers, ensuring that each model generates Chain-of-Thought (CoT) outputs consistent with its intended behavior and style. Since our analysis focuses on the model logits rather than the sampled textual outputs, variations in sampling-related parameters (e.g., temperature, top- k , or top- p) do not affect the logits-based evaluations. This design ensures that our comparisons reflect the models’ intrinsic preference distributions rather than stochastic differences introduced by the decoding process.

All experiments were conducted on two NVIDIA A100-SXM4-80GB GPUs. On average, completing a single experiment—which involves generating step-wise intermediate logits for one LLM on a single dataset using a given Chain-of-Thought (CoT)—takes approximately 20 hours. This reflects the computational demands of step-wise evaluation across multiple inference steps and highlights the resource-intensive nature of detailed logit-level analyses for large language models.

B.2 MCQA Format

The conversion to the MCQA format is illustrated in Table 4. Since MNLI and SNLI belong to the same category of NLI datasets, they are transformed using an identical three-way multiple-choice formulation. In contrast, α NLI is converted using a separate binary-choice MCQA format, reflecting its distinct label structure.

B.3 First-Token-Probability and Metrics

B.3.1 First-token Probability

Take MNLI as an example. Conditioned on the prompts described above, we further map LLM outputs from discrete options in $[A, B, C]$ to probability distributions, which we treat as model judgment distributions (MJDs). Specifically, we define a one-to-one mapping $f: O \rightarrow L$ from the option set O to the label space L , where $O = \{A, B, C\}$ and $L = \{\text{ENTAILMENT}, \text{NEUTRAL}, \text{CONTRADICTION}\}$. Both O and L are subject to permutation to mitigate positional and label-order biases.

Let the textual output of an LLM be represented as a sequence of tokens $w = [w_1, w_2, \dots, w_k]$, where $w_i \in V$, k denotes the output length, and V is the model vocabulary. Instead of using the

| Datasets | MCQA Transformation |
|--------------|---|
| MNLI & SNLI | Please determine whether the following statement is true (entailment), undetermined (neutral), or false (contradiction) given the context below and select ONE of the listed options and start your answer with a single letter. Context: {premise} Statement: {hypothesis} A. Entailment B. Neutral C. Contradiction Answer: |
| α NLI | Please determine which of the two hypotheses (A or B) is more likely to explain the transition from the beginning observation to the ending observation and select ONE of the listed options and start your answer with a single letter. Beginning: {beginning-observation} Ending: {ending-observation} A. {hypothesis1} B. {hypothesis2} Answer: |

Table 4: The MCQA transformation for NLI tasks.

686 decoded output, we extract the pre-decoding logits
687 corresponding to the first generated token w_1 :

$$688 \quad \mathbf{s}_{w_1} = [s_1, s_2, \dots, s_n], \quad n = |V|,$$

689 where s_j denotes the logit associated with the j -th
690 vocabulary token.

691 We restrict our attention to the subset of logits
692 corresponding to the option tokens in O ,

$$693 \quad \mathbf{s}_{w_1}^O = [s_A, s_B, s_C],$$

694 which encode the model’s relative preference over
695 the candidate options. Since the normalization
696 transformation preserves the entropy of the original
697 logits, whereas the softmax transformation (espe-
698 cially when applied with a temperature parame-
699 ter) can alter entropy, we adopt the normalization
700 transformation for our evaluations, because entropy
701 plays a critical role in the computation of JSD, and
702 using a transformation that artificially modifies it
703 could bias the assessment. Therefore, to more ac-
704 curately measure the LLMs’ intrinsic probabilistic
705 preferences and their native reasoning behavior, we
706 rely on the norm transformation rather than softmax
707 in our analysis.

708 To convert these scores into a probability distri-
709 bution p^O , we then apply a normalization step.

$$710 \quad p_{\text{norm}}^O(j) = \frac{s_j}{\sum_{j=1}^{|O|} s_j}, \quad (1)$$

This procedure yields a well-formed probability
distribution over labels, enabling fine-grained com-
parison with human-annotated label distributions.

714 B.3.2 Rank Correlation Metric

715 To quantify the agreement between ranked prefer-
716 ences from different sources (e.g., human annota-
717 tions versus model predictions), we employ rank
718 correlation metrics. Let $\{(x_i, y_i)\}_{i=1}^n$ denote paired
719 ranks from two sources.

720 **Spearman’s ρ** (Spearman, 1961) Spearman’s
721 rank correlation coefficient measures the Pearson
722 correlation between ranked variables and is defined
723 as:

$$724 \quad \rho = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}, \quad (2)$$

725 where $d_i = x_i - y_i$ denotes the rank difference for
726 the i -th item. Spearman’s ρ captures monotonic
727 relationships and is robust to nonlinear transforma-
728 tions of the underlying scores.

729 B.3.3 Distribution-Based Metric

730 For settings where both human annotations and
731 model outputs are represented as probability distri-
732 butions, we adopt distributional similarity metrics.

733 **Jensen–Shannon Distance (JSD)** (Endres and
734 Schindelin, 2003) Given two discrete probability
735 distributions P and Q , the Jensen–Shannon Dis-

tance is defined as:

$$D_{\text{JSD}}(P\|Q) = \sqrt{\frac{1}{2}(D_{\text{KL}}(P\|M) + D_{\text{KL}}(Q\|M))}, \quad (3)$$

where $M = \frac{1}{2}(P + Q)$ and $D_{\text{KL}}(\cdot\|\cdot)$ denotes the Kullback–Leibler divergence. JSD is symmetric, bounded, and well-defined even when P and Q contain zero-probability entries, making it suitable for comparing soft label distributions.

C Box-plot for Cross-CoT Experiments

This section presents the distribution of the Delta JSD metric over all instances in the dataset, aiming to show that the observed reduction in JSD reflects a global trend rather than being driven by a small number of extreme cases that artificially lower the mean.

As illustrated in Figure 4, we visualize per-instance Delta JSD values using box plots, where the central line denotes the median, the boxes correspond to the interquartile range (IQR), and the whiskers extend to $1.5 \times \text{IQR}$. Across the majority of experimental settings, the distributions are centered below zero, indicating a consistent overall decrease in JSD. These results suggest that the improvement captured by the average JSD is broadly shared across data points, rather than being dominated by a few outliers.

D ANOVA Details

To quantify the relative influence of different factors on the observed scores, we employ a two-way Analysis of Variance (ANOVA) with an additive (no-interaction) design. This statistical framework allows us to decompose the total variance of the dependent variable into contributions from multiple categorical factors.

D.1 Problem Setup

Let y_{ij} denote the observed score associated with the i -th model configuration and the j -th CoT setting. In our implementation, we consider:

- A **model factor** with $I = 7$ levels (indexed by i),
- A **CoT factor** with $J = 7$ levels (indexed by j),
- One observation for each (i, j) combination.

The data are organized into a long-form table with three columns: `score`, `model`, and `param`, where both `model` and `param` are treated as categorical variables.

D.2 Additive Two-Way ANOVA Model

We adopt an additive two-way ANOVA model without interaction terms, formulated as:

$$y_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}, \quad (4)$$

where:

- μ is the overall mean score,
- α_i represents the main effect of the i -th model,
- β_j represents the main effect of the j -th CoT,
- ε_{ij} is the residual error term.

This formulation assumes that the effects of the two factors are independent and additive, and that no interaction between model and CoT is modeled. This choice aligns with our goal of isolating the marginal contribution of each factor.

D.3 Estimation via Ordinary Least Squares

The model is estimated using Ordinary Least Squares (OLS), implemented as:

$$\text{score} \sim C(\text{model}) + C(\text{param}), \quad (5)$$

where $C(\cdot)$ denotes categorical encoding. The fitted model is then passed to a Type-II ANOVA procedure, which computes sums of squares for each main effect after accounting for the other factor.

D.4 Variance Decomposition

ANOVA decomposes the total sum of squares (SS) as:

$$SS_{\text{total}} = SS_{\text{model}} + SS_{\text{param}} + SS_{\text{residual}}, \quad (6)$$

where:

- SS_{model} captures variance explained by the model factor,
- SS_{param} captures variance explained by the CoT factor,
- SS_{residual} captures unexplained variance.

Each sum of squares is associated with an F -statistic and corresponding p -value, allowing statistical significance testing of factor effects.

817 D.5 Contribution Percentage

818 To improve interpretability, we further compute the
819 **variance contribution percentage** of each factor:

$$820 \text{Contribution}_k = \frac{SS_k}{SS_{\text{total}}} \times 100\%, \quad (7)$$

821 where $k \in \{\text{model, param, residual}\}$.

822 This metric reflects the proportion of total vari-
823 ance attributable to each source. In our implementa-
824 tion, these percentages are rounded to one decimal
825 place and reported as the final output of the analy-
826 sis.

827 D.6 Interpretation

828 The resulting contribution percentages provide a
829 clear quantitative comparison of how much vari-
830 ability in the scores is explained by:

- 831 • differences between models,
- 832 • differences between CoT settings,
- 833 • unexplained residual noise.

834 This two-way additive ANOVA thus serves as an
835 effective tool for disentangling and comparing the
836 marginal effects of multiple experimental factors
837 in our evaluation framework.

838 E Implementation for Early Stopping

839 E.1 Accumulative 10% Segmenting of 840 Chain-of-Thoughts

841 To facilitate analysis of reasoning progression in
842 Chain-of-Thought (CoT) outputs, we segment each
843 text into ten accumulative portions, corresponding
844 approximately to every 10% of the text length. For-
845 mally, given a text T of length L , we aim to identify
846 cut points p_1, p_2, \dots, p_9 such that p_i roughly cor-
847 responds to $i \cdot L/10$, with the final point $p_{10} = L$.

848 Our procedure combines sentence-aware and
849 heuristic splitting strategies. First, we parse T into
850 sentences using a syntactic parser and extract all
851 sentence-ending positions. If at least nine sentences
852 are available, we select each cut point p_i as the near-
853 est sentence end to $i \cdot L/10$, ensuring monotonicity
854 of cut points. If fewer than ten sentences exist, we
855 iteratively split the longest existing segment, pri-
856 oritizing natural boundaries such as punctuation
857 (e.g., semicolons, commas) and spaces, and resort-
858 ing to midpoint splits when no suitable boundary
859 is found.

860 Finally, we enforce strict monotonicity of cut
861 points and construct the accumulative segments

862 S_1, S_2, \dots, S_{10} where $S_i = T[: p_i]$. This ensures
863 that the last segment always reproduces the full
864 original text. This method preserves original spac-
865 ing and punctuation, providing natural and inter-
866 pretable checkpoints for CoT analysis at decile
867 intervals.

868 E.2 Early-Stopping and Answer Token 869 Extraction

870 To reliably extract intermediate reasoning outputs,
871 we adopt an *early-stopping* strategy. Specifically,
872 at each accumulative CoT segment cut point, we ap-
873 pend a special token sequence signaling the model
874 to terminate reasoning and produce an answer, e.g.,
875 “\n</think>\n\nBased on the reasoning so
876 far, the Answer is:”. This encourages the
877 model to emit the Answer token at that interme-
878 diate stage, preventing incomplete or excessively
879 long continuations.

880 After obtaining the logits for the first token fol-
881 lowing this prompt, we convert them into a proba-
882 bility distribution using the *first-token probability*
883 method. This approach allows us to quantify the
884 model’s intermediate answer distribution at each
885 10% reasoning checkpoint, providing a fine-grained
886 view of decision-making evolution along the CoT.

887 F All Heatmaps for Step-wise Evaluation

888 This section presents the full set of heatmaps ob-
889 tained from our step-wise evaluation, which are
890 subsequently used to compute the ANOVA effect
891 sizes reported in Table 3. Specifically, we provide
892 heatmaps for **accuracy** (Figure 5), **JSD** (Figure 6),
893 **Spearman’s ρ** (Figure 7).

894 G All Curves for Step-wise Evaluation

895 This section presents all curves obtained from the
896 step-wise evaluation, providing a complementary
897 perspective to the heatmaps. Specifically, we show
898 the progression of **accuracy** (Figure 8), **JSD** 9,
899 and **Spearman’s ρ** (Figure 10) across inference
900 steps. These curves allow us to track how each
901 metric evolves throughout the reasoning process,
902 revealing dynamic trends in model performance,
903 distributional alignment, and rank correlation over
904 time.

905 H Analysis of CoT Formats

906 To further investigate the “split influence” observed
907 in our quantitative results—where CoT determines

908 the final accuracy but leaves the distributional struc-
909 ture (JSD and Spearman’s ρ) largely anchored to
910 model priors—we conducted an examination of the
911 generated reasoning traces.

912 As discussed in Section 4.2, step-wise analysis
913 reveals that accuracy often shifts or converges only
914 at the conclusion, while Spearman’s ρ fluctuates
915 without a consistent trend. Thus, CoT determines
916 the LLM’s final choice but not the structure of un-
917 certainty. We attribute this to the structural format
918 of CoT: reasoning traces typically culminate in ex-
919 plicit, decisive conclusion statements (e.g., “There-
920 fore, the answer is...”), which strongly steer the
921 accuracy in the final steps.

922 Table 5 presents anecdotal evidence from the
923 α NLI dataset. We observe a consistent pattern
924 across models:

- 925 • **Explicit Conclusion at the End:** The reason-
926 ing process consistently ends with a strong,
927 definitive statement identifying the correct op-
928 tion (e.g., “So, A must be the correct choice”).
929 This explicit signal aligns with the sharp rise
930 in accuracy observed in the final steps of our
931 step-wise analysis.
- 932 • **Implicit Distributional Weighing:** While the
933 models argue *for* the best option, they rarely
934 explicitly articulate the relative probability of
935 the alternative options in a way that would
936 restructure the output distribution. Conse-
937 quently, while the final decision is explicitly
938 dictated by the CoT’s conclusion, the distri-
939 butional structure over non-argmax options
940 remains latent and implicit, governed largely
941 by the model’s intrinsic priors.

| Model & Case | Reasoning Excerpt (Conclusion Phase) | Observation regarding Final Choice vs. Structure |
|-------------------------------------|---|---|
| Qwen (Instance 2: Sandy) | “...Therefore, A is better. I think B is a distractor. [...] The instruction says ‘select ONE of the listed options’ ... So, my response should be just ‘A’. ” | Decisive Locking: The reasoning concludes by explicitly discarding the alternative and locking onto the single target token ‘A’, driving the final accuracy without refining the relative probability space. |
| GLM (Instance 3: Bananas) | “...Therefore, B cannot explain why they’re talking about eating a banana. So, A must be the correct choice. [...] Therefore, A is correct. ” | Convergence to Argmax: The trace culminates in strong assertions (“must be”, “correct”), which serve to fix the model’s final decision, explaining why accuracy converges at the end while latent distributions remain implicit. |
| GPT (Instance 1: Ron) | “...That suggests his actions. So A is more appropriate. Thus choose A. [...] We must start answer with a single letter... So final answer: ‘A’. ” | Explicit Selection: The reasoning shifts from semantic evaluation to an operational selection command (“Thus choose A”), confirming that the CoT acts as a decision-maker for the top option. |

Table 5: Examples of CoT reasoning traces from the α NLI dataset. The excerpts illustrate how CoT reasoning typically ends with explicit conclusion statements. This structural characteristic supports our finding that CoT content determines the final choice (Accuracy) through explicit reasoning, while the underlying structure of uncertainty (JSD/Ranking) remains latent and less affected by these definitive concluding remarks.

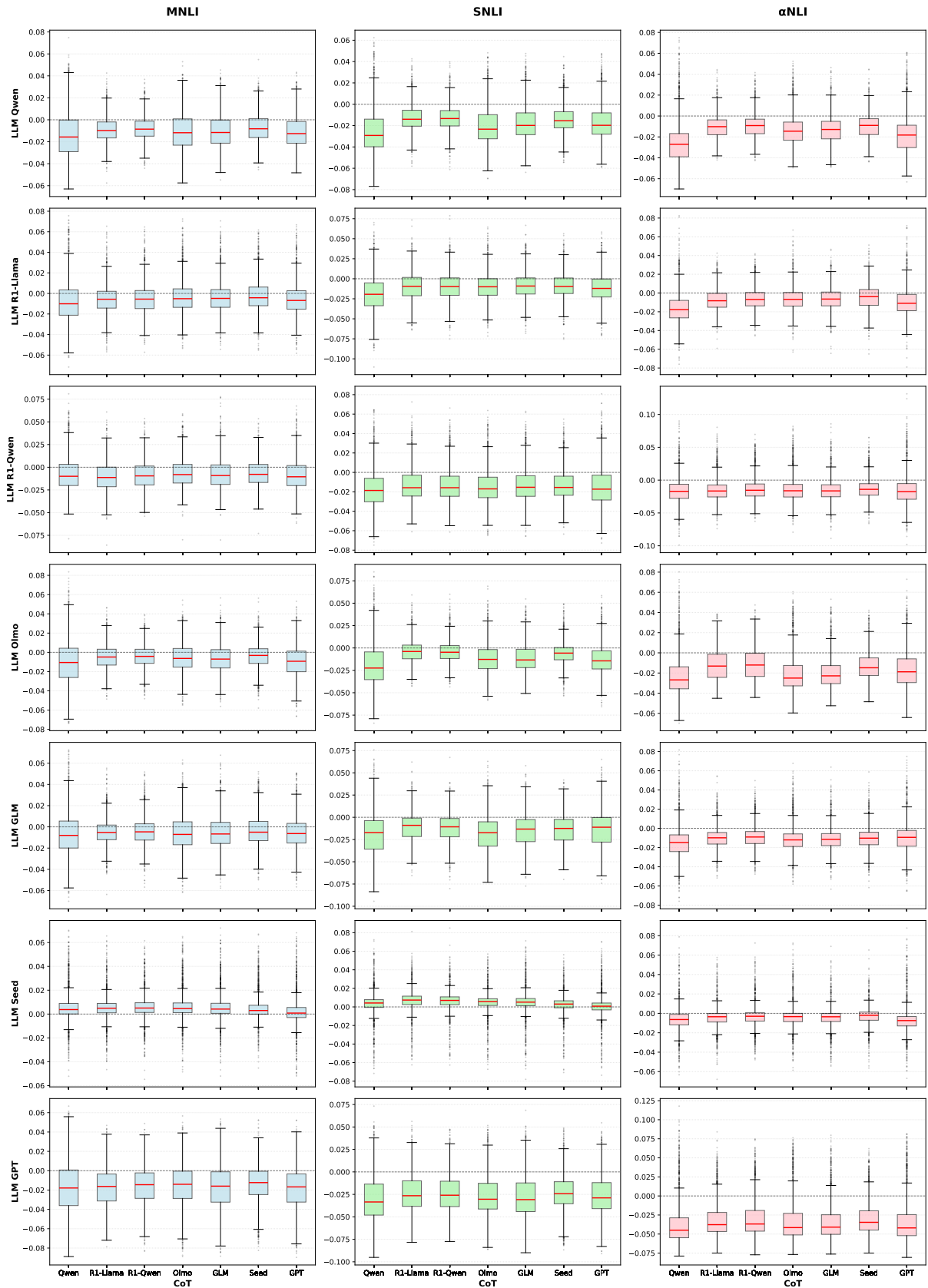


Figure 4: Delta JSD box plot.

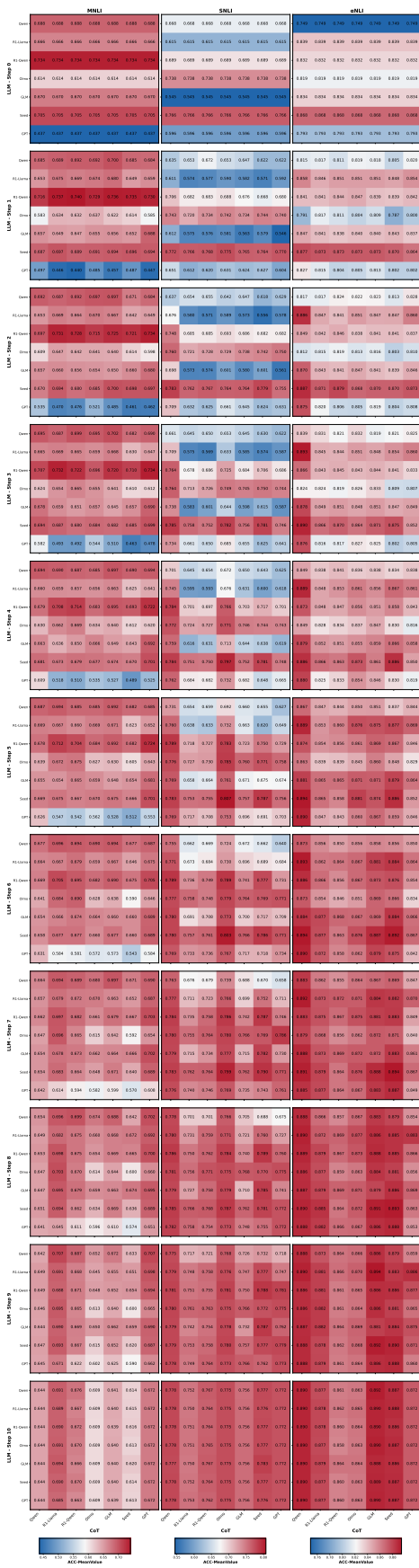


Figure 5: Steps ACC.



Figure 6: Steps JSD.

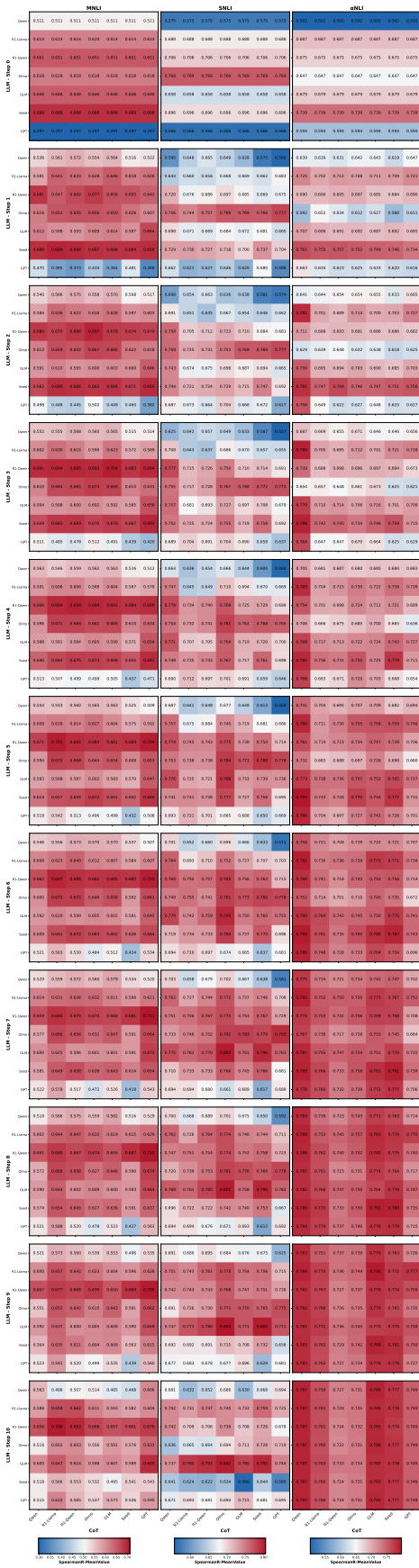


Figure 7: Steps Spearman's ρ .

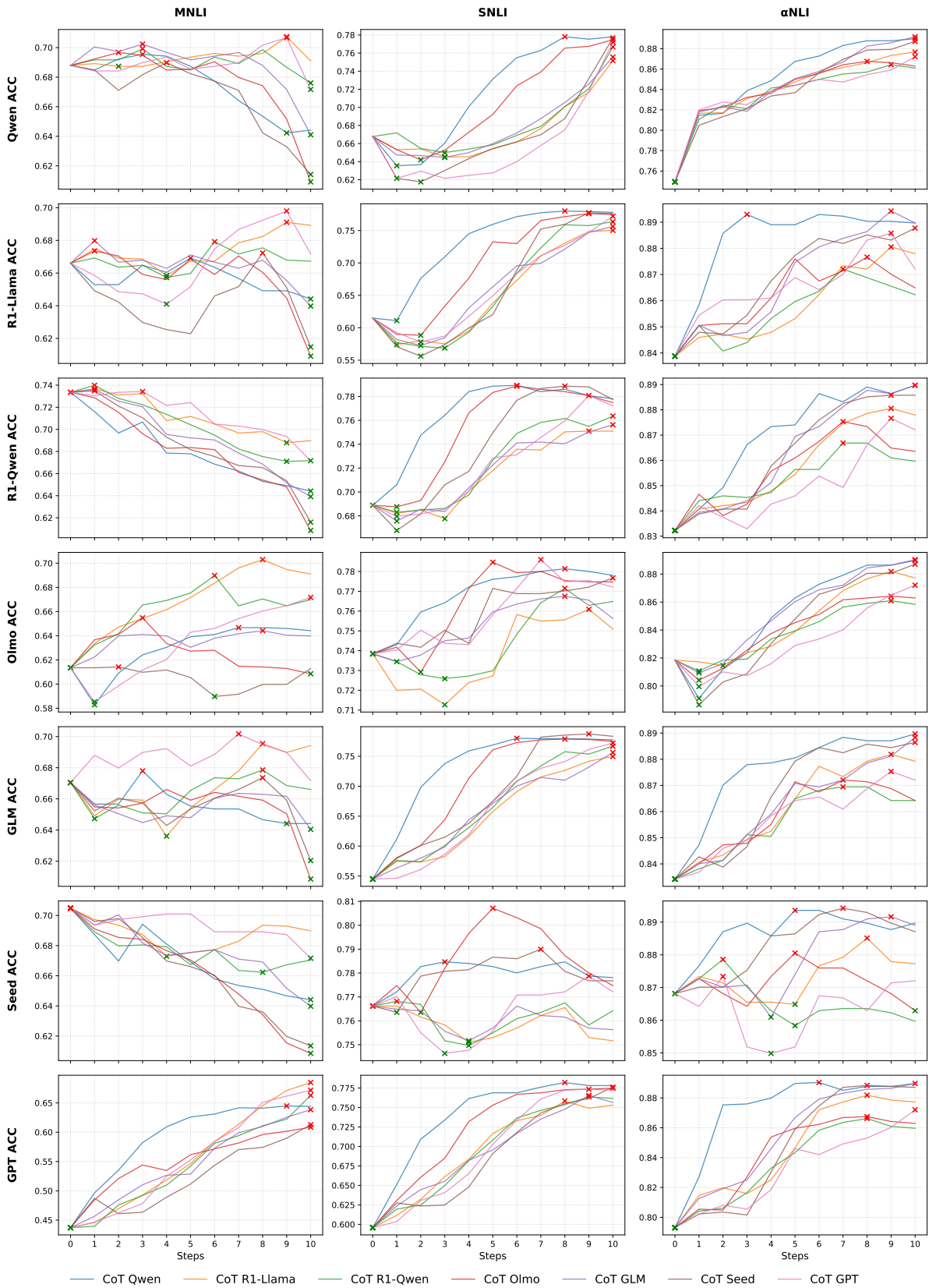


Figure 8: Curves ACC.

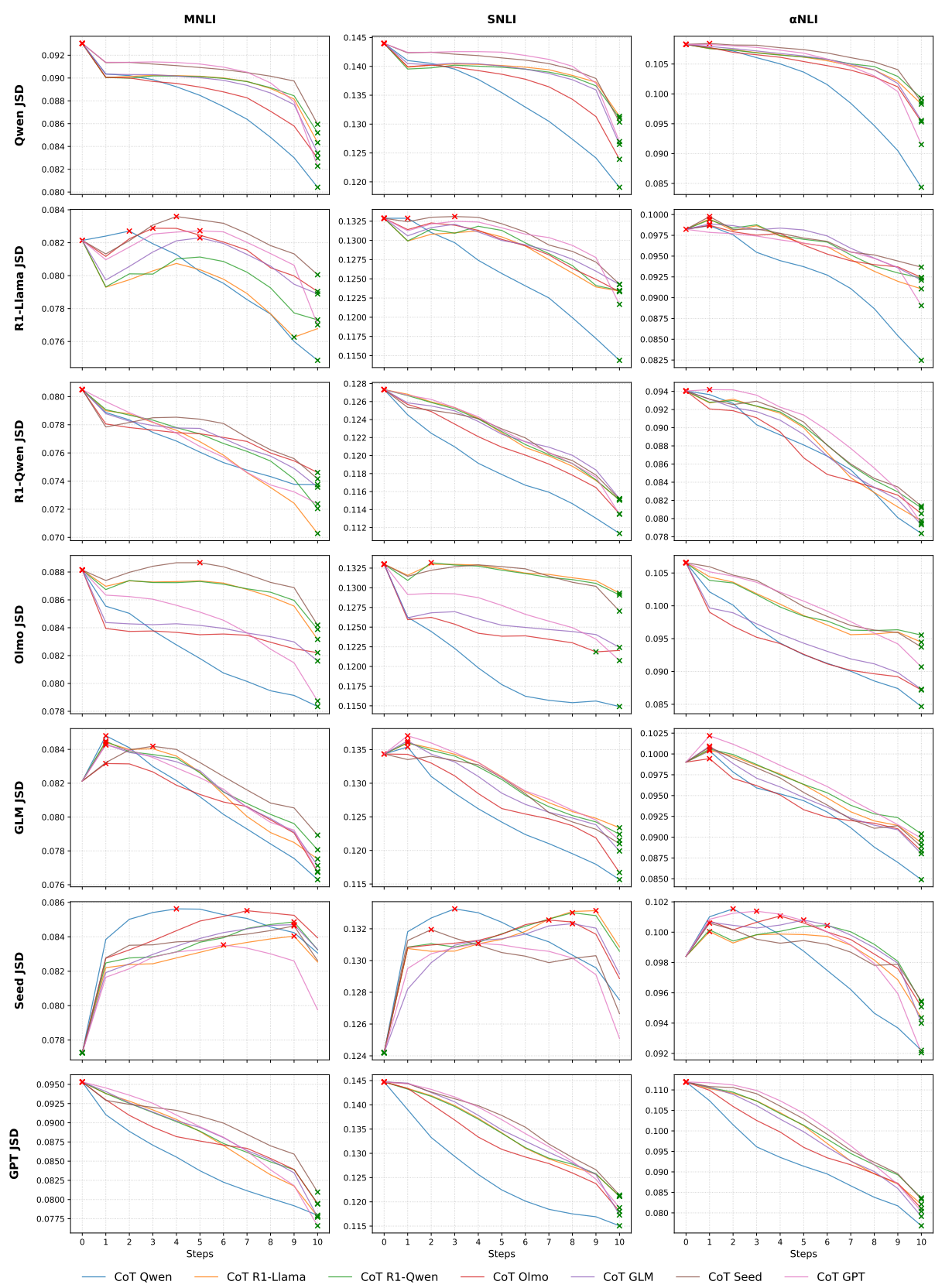


Figure 9: Curves JSD.

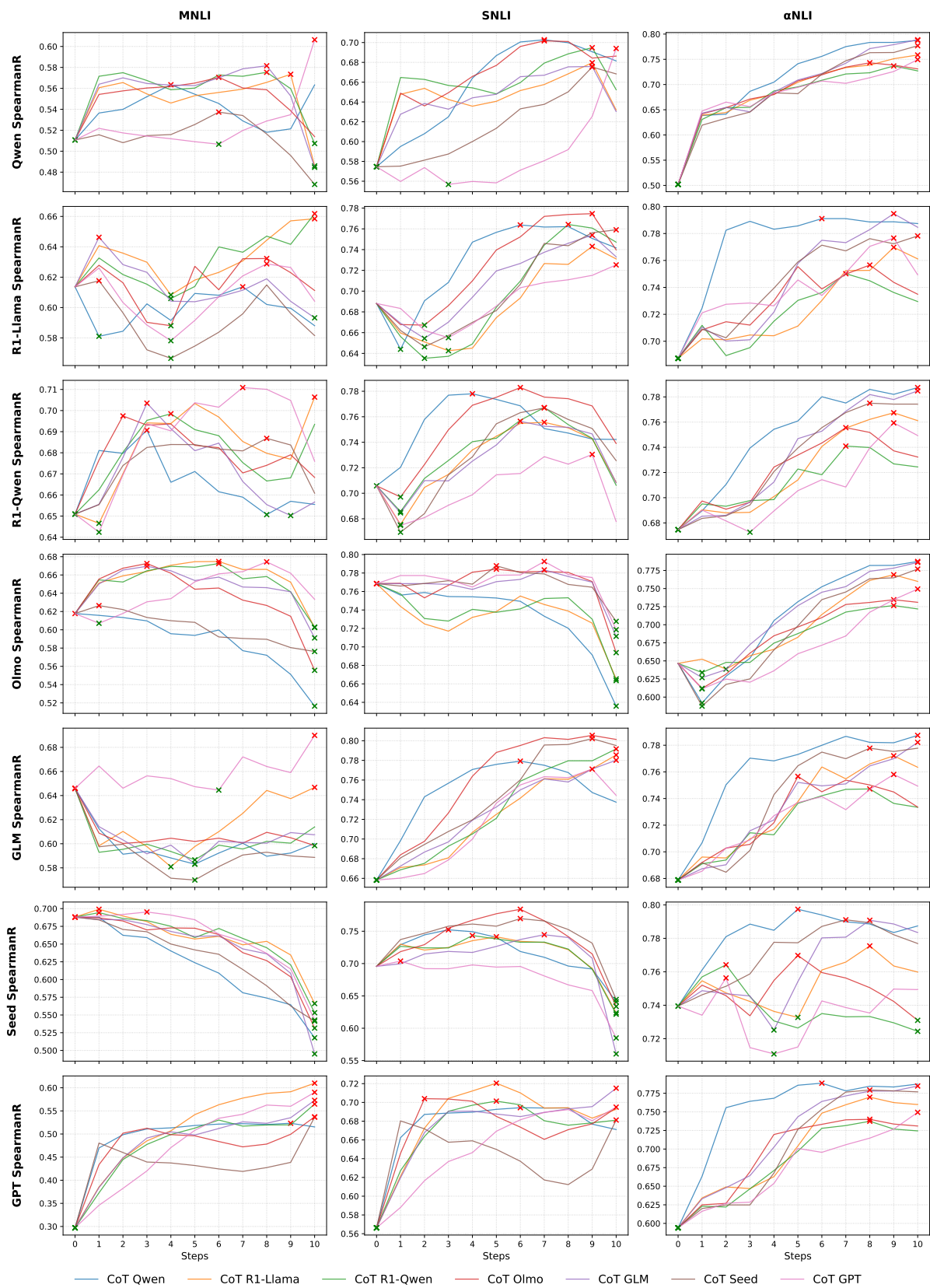


Figure 10: Curves Spearman's ρ .