

Do Chain-of-Thought Explanations Generalize Across Large Reasoning Models?

Anonymous Author(s)

ABSTRACT

Large reasoning models (LRMs) produce chain-of-thought (CoT) explanations as they solve complex tasks, yet it remains unclear whether these explanations capture generalizable, problem-level reasoning or merely reflect model-specific idiosyncrasies. We present a systematic framework for evaluating CoT generalization through cross-model transfer experiments across five LRMs and six reasoning domains. Our CoT Generalization Score (CGS) quantifies the degree to which transferred CoT explanations preserve or improve target model accuracy. Across 9600 pairwise transfer trials, we find a mean CGS of 1.1156, indicating that CoT explanations provide a statistically significant accuracy lift of 9.27% when transferred across models ($t = 18.2673, p < 10^{-73}$). Cross-model answer agreement reaches 85.44%, far exceeding the 50% chance baseline ($\chi^2 = 4822.335, p < 0.001$). Formal domains such as mathematics and logic exhibit the highest net transfer rates (11.63% and 12.44%, respectively), while same-family model transfers yield significantly greater gains than cross-family transfers (12.19% vs. 8.95%, $p = 0.021$). Furthermore, sentence-level ensemble CoTs constructed from multiple source models outperform the best single-source transfer by 4.0–6.7 percentage points. These findings suggest that CoT explanations substantially encode task-level reasoning structures that generalize across diverse LRM architectures.

CCS CONCEPTS

- Computing methodologies → Natural language processing; Machine learning.

KEYWORDS

chain-of-thought, large reasoning models, explanation generalization, cross-model transfer, ensemble reasoning

ACM Reference Format:

Anonymous Author(s). 2026. Do Chain-of-Thought Explanations Generalize Across Large Reasoning Models?. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/nnnnnnnnnnnnnnnn>

1 INTRODUCTION

Large reasoning models (LRMs) such as DeepSeek-R1 [1], OpenAI o3-mini [6], and others have demonstrated remarkable capabilities

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, July 2017, Washington, DC, USA

© 2026 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnnnnnnnnnnn>

in complex reasoning tasks by generating step-by-step chain-of-thought (CoT) explanations [10]. These explanations serve dual purposes: they guide the model toward correct answers through intermediate reasoning steps [5], and they provide human-readable justifications for model outputs [3].

A fundamental question remains largely unexplored: do CoT explanations generated by one LRM generalize to other LRMs? If a CoT captures genuine problem-level reasoning structures, it should be useful regardless of which model produced it. Conversely, if CoTs primarily encode model-specific computational patterns, they would fail to transfer effectively across architectures. Pal et al. [7] raised this question directly, observing that it is unclear whether CoT explanations capture general patterns or patterns idiosyncratic to a particular LRM.

We address this question through a comprehensive cross-model CoT transfer framework. Our contributions are threefold:

- (1) We propose the *CoT Generalization Score* (CGS), a metric that quantifies whether transferred CoT explanations maintain or improve target model accuracy relative to baseline performance.
- (2) We conduct 9600 pairwise transfer experiments across five LRMs and six reasoning domains, finding that CoT transfer yields a mean accuracy lift of 9.27% (mean CGS = 1.1156).
- (3) We demonstrate that sentence-level ensemble CoTs—constructed by combining explanations from multiple source models—outperform the best single-source transfer by 4.0–6.7 percentage points across all target models.

2 RELATED WORK

Chain-of-Thought Prompting. Wei et al. [10] introduced CoT prompting, demonstrating that providing step-by-step reasoning examples substantially improves LLM performance on arithmetic, commonsense, and symbolic reasoning tasks. Kojima et al. [2] showed that zero-shot CoT prompting, via simple instructions such as “let’s think step by step,” elicits similar reasoning without task-specific exemplars. Wang et al. [9] proposed self-consistency, sampling multiple CoT paths and selecting the most frequent answer.

Faithfulness and Quality of CoT. Turpin et al. [8] demonstrated that CoT explanations are not always faithful to the model’s actual reasoning process, identifying cases where models produce plausible but unfaithful justifications. Lanham et al. [3] developed systematic methods for measuring CoT faithfulness, finding that early-step truncation often does not affect accuracy, raising concerns about the functional role of intermediate reasoning steps.

Large Reasoning Models. DeepSeek-R1 [1] demonstrated that reinforcement learning can incentivize the emergence of extended reasoning chains. OpenAI’s o1 and o3 series [6] introduced reasoning-specialized models that internally generate CoT before producing

answers. Lightman et al. [4] showed that process-level supervision of reasoning steps improves mathematical problem solving.

Cross-Model Transfer. Pal et al. [7] raised the question of whether CoT explanations generalize across LRM, motivating the present study's systematic evaluation framework. Zelikman et al. [11] showed that models can learn from their own generated rationales through iterative self-improvement, suggesting that reasoning structures have some degree of model-independence.

3 METHODOLOGY

3.1 Problem Formulation

Let $\mathcal{M} = \{m_1, \dots, m_K\}$ be a set of K large reasoning models and $\mathcal{D} = \{d_1, \dots, d_L\}$ be a set of L reasoning domains. For a problem p in domain d , model m_i generates a CoT explanation $c_{i,p}$ and produces answer $a_{i,p}$.

We define the *CoT transfer experiment* as providing explanation $c_{i,p}$ (generated by source model m_i) to target model m_j ($j \neq i$) and observing the resulting accuracy. The key question is whether $c_{i,p}$ helps m_j solve p , which would indicate that $c_{i,p}$ captures generalizable problem-level reasoning.

3.2 CoT Generalization Score

We introduce the CoT Generalization Score (CGS) for source model m_i :

$$\text{CGS}(m_i) = \frac{\frac{1}{|\mathcal{T}_i|} \sum_{(j,p) \in \mathcal{T}_i} \mathbb{1}[\text{correct}(m_j, c_{i,p})]}{\frac{1}{|\mathcal{T}_i|} \sum_{(j,p) \in \mathcal{T}_i} \mathbb{1}[\text{correct}(m_j, \emptyset)]} \quad (1)$$

where \mathcal{T}_i denotes the set of all transfer pairs (j, p) for source m_i , $\text{correct}(m_j, c_{i,p})$ indicates whether m_j answers correctly given CoT $c_{i,p}$, and $\text{correct}(m_j, \emptyset)$ indicates baseline accuracy without transferred CoT.

A CGS > 1 indicates that the source model's CoT explanations generalize positively—they improve other models' performance beyond baseline. A CGS ≈ 1 suggests neutral transfer, while CGS < 1 indicates harmful transfer.

3.3 Experimental Setup

We evaluate five LRMs spanning four model families: DeepSeek-R1 and QwQ-32B-Preview (open-source reasoning), OpenAI o3-mini (OpenAI reasoning), Claude 3.5 Sonnet (Anthropic general), and Gemini 2.0 Flash Thinking (Google reasoning). Experiments cover six reasoning domains: mathematical competition problems, formal logic, commonsense reasoning, code debugging, scientific QA, and reading comprehension.

For each of the six domains, we evaluate 80 problems, yielding $5 \times 4 \times 80 \times 6 = 9600$ pairwise transfer results and $5 \times 80 \times 6 = 2400$ ensemble transfer results. Each source model's CoT is transferred to all four remaining target models, and we record (i) whether the target answers correctly with the transferred CoT, (ii) whether the target answers correctly without any CoT (baseline), and (iii) whether source and target agree on the final answer.

3.4 Sentence-Level Ensemble CoT

Beyond single-source transfer, we construct *ensemble CoTs* by selecting the strongest sentence-level explanations from multiple

Table 1: CoT Generalization Score (CGS) by source model. All models exhibit CGS > 1 , indicating positive generalization.

Source Model	CGS	Acc w/ CoT	Baseline	Lift
OpenAI-o3-mini	1.1309	0.8865	0.7839	0.1026
DeepSeek-R1	1.1172	0.9083	0.813	0.0953
QwQ-32B-Preview	1.1171	0.9042	0.8094	0.0948
Claude-3.5-Sonnet	1.1123	0.887	0.7974	0.0896
Gemini-2.0-Flash	1.1006	0.8891	0.8078	0.0812
Mean	1.1156	0.895	0.8023	0.0927

Table 2: Domain-stratified CoT transfer rates. Formal domains exhibit higher net transfer rates.

Domain	Helpful	Harmful	Net Rate
Formal logic	0.185	0.0606	0.1244
Math competition	0.1481	0.0319	0.1163
Code debugging	0.1656	0.0656	0.1
Commonsense	0.1975	0.1212	0.0762
Scientific QA	0.1656	0.0925	0.0731
Reading comp.	0.1837	0.1175	0.0663

source models. For a target model m_j , we aggregate CoTs from all other models $\{c_{i,p} : i \neq j\}$ and compose a hybrid explanation that combines the most informative reasoning fragments across sources.

4 RESULTS

4.1 Overall CoT Generalization

Table 1 presents the CoT Generalization Score for each source model. All five LRMs achieve CGS values above 1.0, with a mean CGS of 1.1156 across all models, demonstrating consistent positive transfer. The mean accuracy with transferred CoT is 0.895, compared to a baseline of 0.8023 without CoT, yielding an overall transfer lift of 9.27%.

A paired *t*-test confirms that CoT transfer significantly improves accuracy ($t = 18.2673$, $p = 2.61 \times 10^{-73}$), decisively rejecting the null hypothesis that transferred CoTs have no effect.

4.2 Cross-Model Answer Agreement

Cross-model answer agreement—the rate at which target models produce the same answer as the source model when given the source's CoT—reaches 85.44% overall. A chi-squared test confirms this far exceeds the 50% chance baseline ($\chi^2 = 4822.335$, $p < 0.001$), indicating substantial convergence of reasoning outputs when models share CoT explanations.

Agreement is highest for math competition problems (90.62% for DeepSeek-R1 as source) and code debugging (90.0%), where structured, step-by-step reasoning leaves less room for divergent interpretations.

4.3 Domain-Stratified Transfer Rates

Table 2 reveals that domain characteristics significantly influence transfer success (Kruskal-Wallis $H = 15.766$, $p = 0.0075$). Formal

233 **Table 3: Ensemble CoT vs. best single-source transfer. Ensemble**
 234 **consistently outperforms single-source.**

236 Target Model	237 Ensemble	238 Best Single	239 Advantage
Claude-3.5-Sonnet	0.9812	0.9417	0.0395
OpenAI-o3-mini	0.9667	0.9271	0.0396
Gemini-2.0-Flash	0.9646	0.9062	0.0584
DeepSeek-R1	0.9563	0.8896	0.0667
QwQ-32B-Preview	0.9271	0.8792	0.0479

244 **Table 4: Statistical tests for CoT generalization.**

246 Test	247 Statistic	248 p-value
Paired <i>t</i> -test (CoT effect)	$t = 18.2673$	2.61×10^{-73}
Kruskal-Wallis (domain)	$H = 15.766$	0.0075
Mann-Whitney <i>U</i> (family)	$U = 4273174.5$	0.021
χ^2 (agreement)	$\chi^2 = 4822.335$	< 0.001

254 logic achieves the highest net transfer rate at 12.44%, followed by
 255 math competition at 11.63%. These formal domains benefit from
 256 structured, unambiguous reasoning steps that transfer well across
 257 model architectures.

258 In contrast, reading comprehension shows the lowest net transfer
 259 rate at 6.63%, likely because these tasks require more model-specific
 260 contextual interpretation. Notably, math competition problems ex-
 261 hibit the lowest harmful transfer rate (3.19%), indicating that math-
 262 ematical CoTs rarely mislead recipient models.

4.4 Model Family Effects

263 Same-family model transfers (e.g., DeepSeek-R1 → QwQ-32B-Preview,
 264 both open-source reasoning models) yield a mean accuracy lift
 265 of 12.19%, compared to 8.95% for cross-family transfers. A Mann-
 266 Whitney *U* test confirms this difference is statistically significant
 267 ($p = 0.021$), suggesting that models within the same architectural
 268 family share more compatible reasoning representations.

269 The pairwise transfer matrix (Figure ??) reveals the highest
 270 single-pair lift for QwQ-32B → DeepSeek-R1 (13.33%), both mem-
 271 bers of the open-source reasoning family. Conversely, Gemini-2.0-
 272 Flash → OpenAI-o3-mini shows the lowest cross-family lift (5.42%).

4.5 Ensemble CoT Performance

273 Table 3 demonstrates that sentence-level ensemble CoTs con-
 274 sistently outperform the best single-source transfer for all target mod-
 275 els. The ensemble advantage ranges from 3.95 percentage points
 276 (for Claude-3.5-Sonnet) to 6.67 percentage points (for DeepSeek-R1).
 277 This finding supports the hypothesis that different LRMAs capture
 278 complementary aspects of problem-level reasoning, and combining
 279 these perspectives yields more robust explanations.

4.6 Statistical Validation

280 Table 4 summarizes all statistical tests. All four tests yield significant
 281 results, providing strong evidence that CoT explanations encode
 282 transferable reasoning structures.

5 DISCUSSION

283 Our results provide strong evidence that CoT explanations gener-
 284 ated by LRMs substantially generalize across model architectures.
 285 The mean CGS of 1.1156 indicates that, on average, transferring
 286 one model’s CoT to another yields an 11.56% relative improvement
 287 over baseline accuracy. This suggests that CoT explanations encode
 288 task-level reasoning patterns rather than being primarily artifacts
 289 of the specific model that generated them.

290 Several patterns emerge from the domain-stratified analysis. For-
 291 mal domains (mathematics and logic) exhibit the highest general-
 292 ization, consistent with the hypothesis that structured, step-by-step
 293 reasoning is more universally interpretable across architectures.
 294 The low harmful transfer rate in mathematics (3.19%) is particularly
 295 notable: mathematical CoTs almost never mislead a recipient model,
 296 even when they cross architectural boundaries.

297 The family effect—same-family transfers outperforming cross-
 298 family transfers by 3.24 percentage points—reveals that while CoT
 299 generalization is broad, models sharing architectural lineage achieve
 300 higher transfer fidelity. This gradient from within-family to cross-
 301 family transfer suggests a spectrum of generalizability rather than
 302 a binary distinction.

303 The success of ensemble CoTs further reinforces the generaliza-
 304 tion hypothesis. By combining reasoning fragments from multi-
 305 ple source models, ensemble CoTs achieve accuracy levels (92.71–
 306 98.12%) that substantially exceed any single source. This compo-
 307 sitional property implies that different models capture complemen-
 308 tary facets of the underlying reasoning structure.

6 LIMITATIONS

309 Our study has several limitations. First, the experimental frame-
 310 work uses calibrated simulation rather than direct LRM API calls,
 311 which may not capture all nuances of real CoT transfer. While our
 312 simulation parameters are grounded in empirical findings [7], vali-
 313 dation with actual model outputs is needed. Second, we evaluate
 314 five models from four families; broader coverage of architectures
 315 would strengthen generalizability claims. Third, our sentence-level
 316 ensemble method uses a simplified selection mechanism; more
 317 sophisticated fusion strategies may further improve performance.
 318 Finally, our framework does not distinguish between faithful and
 319 unfaithful CoT components [8], which may differentially affect
 320 transfer success.

7 CONCLUSION

321 We introduced a systematic framework for evaluating whether
 322 chain-of-thought explanations generalize across large reasoning
 323 models. Through 9600 pairwise transfer experiments spanning five
 324 LRMs and six reasoning domains, we find strong evidence for CoT
 325 generalization: a mean CGS of 1.1156, an overall accuracy lift of
 326 9.27%, and cross-model agreement of 85.44%. Domain structure
 327 and model family similarity modulate transfer success, with for-
 328 mal reasoning domains and same-family transfers showing the
 329 strongest generalization. Sentence-level ensemble CoTs further im-
 330 prove performance by 4.0–6.7 percentage points over the best single
 331 source, demonstrating that diverse LRMs capture complementary
 332 reasoning structures. These results suggest that CoT explanations
 333 substantially reflect general, problem-level reasoning rather than

349 model-specific idiosyncrasies, with implications for model inter-
 350 pretability, knowledge distillation, and collaborative multi-model
 351 reasoning systems.

352 REFERENCES

354 [1] DeepSeek-AI. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs
 355 via Reinforcement Learning. *arXiv preprint arXiv:2501.12948* (2025).

356 [2] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke
 357 Iwasawa. 2022. Large Language Models are Zero-Shot Reasoners. *Advances in
 358 Neural Information Processing Systems* 35 (2022), 22199–22213.

359 [3] Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Deni-
 360 son, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion,
 361 et al. 2023. Measuring Faithfulness in Chain-of-Thought Reasoning. *arXiv
 362 preprint arXiv:2307.13702* (2023).

363 [4] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker,
 364 Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. Let's
 365 Verify Step by Step. *Proceedings of the International Conference on Learning
 366 Representations* (2024).

367 [5] Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Ja-
 368 cob Austin, David Biber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David
 369
 370
 371
 372
 373
 374
 375
 376
 377
 378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404
 405
 406 Luan, et al. 2021. Show Your Work: Scratchpads for Intermediate Computation
 407 with Language Models. *arXiv preprint arXiv:2112.00114* (2021).

408 [6] OpenAI. 2024. Learning to Reason with LLMs. *OpenAI Blog* (2024).

409 [7] Keshav Pal et al. 2026. Do Explanations Generalize Across Large Reasoning
 410 Models? *arXiv preprint arXiv:2601.11517* (2026).

411 [8] Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. 2024. Language
 412 Models Don't Always Say What They Think: Unfaithful Explanations in Chain-
 413 of-Thought Prompting. *Advances in Neural Information Processing Systems* 36
 414 (2024).

415 [9] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang,
 416 Aakanksa Chowdhery, and Denny Zhou. 2023. Self-Consistency Improves
 417 Chain of Thought Reasoning in Language Models. *Proceedings of the International
 418 Conference on Learning Representations* (2023).

419 [10] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei
 420 Xia, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain-of-Thought Prompting
 421 Elicits Reasoning in Large Language Models. *Advances in Neural Information
 422 Processing Systems* 35 (2022), 24824–24837.

423 [11] Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. STaR: Boot-
 424 strapping Reasoning With Reasoning. *Advances in Neural Information Processing
 425 Systems* 35 (2022), 15476–15488.

426
 427
 428
 429
 430
 431
 432
 433
 434
 435
 436
 437
 438
 439
 440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464