

A Survey on Chain-of-Thought Reasoning Evaluation in Large Language Models: What, How, Who, and Where

Anonymous ACL submission

Abstract

As large language models evolve from conversational agents to reasoning engines, Chain-of-Thought (CoT) evaluation has become pivotal, yet remains fragmented and heavily reliant on final accuracy. This outcome bias often masks *disguised accuracy*, failing to distinguish faithful reasoning from post-hoc rationalization. Despite the urgency, a comprehensive survey systematizing these process-oriented assessment techniques remains absent. To fill this gap, we present a unified framework organizing the literature along four dimensions: *what to evaluate*, *how to evaluate*, *who evaluates*, and *where to evaluate*. We formalize core quality metrics and systematically analyze evaluation methodologies across qualitative and quantitative paradigms within diverse application domains. By synthesizing these approaches to improve CoT reasoning reliability, this survey provides guidance for building robust evaluation pipelines and highlights key frontiers for future research.

1 Introduction

Recently, large language models (LLMs), such as GPT-5 (OpenAI, 2025) and Gemini-3 (Google DeepMind, 2025), have advanced toward sophisticated Chain-of-Thought (CoT) reasoning (Wei et al., 2022; Yang et al., 2025a; Guo et al., 2025a). This paradigm alleviates the traditional black-box nature of LLMs (Bommasani et al., 2021; Rae et al., 2021) by explicitly articulating intermediate reasoning steps akin to human cognition (Chu et al., 2024; Chen et al., 2025c). Such transparency not only enhances interpretability (Ye et al., 2023b; Chen et al., 2024a) but also improves performance across diverse tasks (Cobbe et al., 2021; Kojima et al., 2022; Zhang et al., 2023b; Nye et al., 2021; Wang et al., 2023c; Talmor et al., 2023).

However, as shown in Figure 1, this shift toward CoT reasoning introduces new challenges in assessing the reliability of the reasoning process

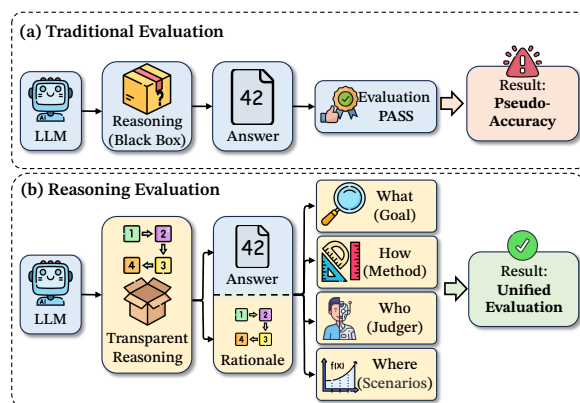


Figure 1: Unlike traditional outcome-based surveys, our survey systematically summarizes transparent CoT reasoning evaluation methodologies.

itself, beyond merely evaluating the correctness of final outputs (Madaan et al., 2023b; Lee and Hockenmaier, 2025; Zhao et al., 2025b). Conventional evaluation metrics, which primarily emphasize answer accuracy (Rajpurkar et al., 2016; Wang et al., 2018), fail to capture the nuanced quality of intermediate steps that underpin model behavior. Crucially, this ignores *disguised accuracy*¹, where models reach correct answers via flawed logic. Such outcome-bias leads to evaluations that fail to reflect the model’s true capabilities.

Consequently, establishing reliable CoT evaluation is a central objective (Lee and Hockenmaier, 2025), essential for verifying reasoning soundness behind correct predictions and systematically characterizing LLM capabilities. Recent studies propose diverse approaches, from step-level assessments of process fidelity (Lightman et al., 2024; Uesato et al., 2022; Zhang et al., 2025d) to model-based judgment frameworks (Zheng et al., 2023; Kim et al., 2024b; Leang et al., 2025), to quantify such capabilities (Guan et al., 2025; Guo et al., 2025b). Despite these advances, current efforts

¹An example is shown in Appendix A.

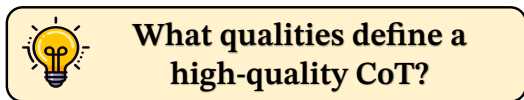
066 remain conceptually fragmented, and no compre- 111
067 hensive survey has yet unified evaluation principles 112
068 across existing methodologies.

069 To fill this blank, we present the first compre- 113
070 hensive survey of CoT evaluation, organized under an 114
071 integrative *What-How-Who-Where* structure. We 115
072 begin by defining *what to evaluate* (§2), extending 116
073 the notion of trustworthy reasoning beyond surface- 117
074 level accuracy to encompass multiple cognitive 118
075 and behavioral dimensions. We then classify *how* 119
076 *to evaluate* (§3) across qualitative evaluation and 120
077 quantitative evaluation, highlighting their method- 121
078 ological trade-offs and empirical foundations. Next, 122
079 we examine *who evaluates* (§4), assessing the bal- 123
080 ance between human judgment and automated scal- 124
081 ability in ensuring evaluation reliability. Finally, 125
082 we map *where to evaluate* (§5) to application do- 126
083 mains ranging from closed-world logical inference 127
084 to open-ended dialogue and scientific reasoning. 128
085 The concise taxonomy underpinning this frame- 129
086 work is illustrated in Figure 2. Furthermore, we 130
087 discuss emerging challenges and future directions 131
088 (§6), aiming to guide the development of reliable 132
089 CoT evaluation methodologies. 133

090 Our contributions are summarized as follows:

- 091 • **First Survey:** We present the first compre- 134
092 hensive survey on evaluating CoT reasoning. 135
- 093 • **Systematic Taxonomy:** We critically analyze 136
094 current evaluation methodologies, identifying 137
095 What-How-Who-Where dimensions in assess- 138
096 ing the reliability of CoT reasoning. 139
- 097 • **New Frontiers:** We outline future research di- 140
098 rections to advance the field of CoT evaluation, 141
099 emphasizing the need for standardized bench- 142
100 marks and more nuanced assessment criteria. 143

101 2 What to Evaluate 144



102 Transcending outcome-based metrics that overlook 145
103 reasoning nuances, Figure 3 illustrates our unified 146
104 framework structured along five core dimensions. 147

105 2.1 Correctness 148

106 Correctness refers to the overall validity of CoT, 149
107 capturing whether the final conclusion or result is 150
108 factually and computationally sound (Wei et al., 151
109 2022; Kojima et al., 2022). It represents a founda- 152
110 tional criterion for reasoning quality, as correctness 153

111 establishes the necessary condition for any valid 112
113 outcome (Cobbe et al., 2021; Uesato et al., 2022).

114 Formally, let a reasoning chain be defined as 115
116 an ordered sequence of inferential steps $R =$ 117
118 $\{r_1, r_2, \dots, r_n\}$, culminating in a final answer a . 119
120 The reasoning chain satisfies overall correctness if 121
122 and only if $\phi(a) = \text{True}$, where ϕ denotes a val- 123
124 idation function that ensures factual, logical, and 125
126 computational integrity. 127

120 2.2 Coherence 120

121 Coherence assesses logic structural soundness, en- 122
123 suring deductive validity in reasoning irrespective 124
125 of factual accuracy (Pan et al., 2023; Golovneva 126
127 et al., 2023). Generally speaking, flawed logic 128
129 can yield hallucinated conclusions (Zhang et al., 130
131 2024b), necessitating step-wise fidelity over mere 132
133 overall correctness (Madaan et al., 2023b). Process 134
135 reward models detect local non-sequiturs in verified 136
137 chains (Lightman et al., 2024; Uesato et al., 2022). 138
139 Symbolic solvers and logic checkers provide deter- 140
141 minism and precision (Lyu et al., 2023a; Zhao et al., 142
143 2025b). Meanwhile, natural language inference 144
145 models are utilized to measure step entailment and 146
147 uncover cycles or disconnected reasoning (Zhao 148
149 et al., 2025b; Valmeekam et al., 2022). Coherence 150
151 thus anchors rational deduction, distinguishing it 152
153 from narrative association. 154

155 For the reasoning steps $(r_i, r_{i+1}) \in R$, coher- 156
157 ence is satisfied when entailment $r_i \models r_{i+1}$ holds 158
159 for all i . Formally, logical coherence is defined as:

$$160 \mathcal{C}(R) = \frac{1}{n-1} \sum_{i=1}^{n-1} \text{Entail}(r_i, r_{i+1}), \quad (1)$$

161 where $\text{Entail}(\cdot)$ returns 1 if reasoning is valid. 162

163 2.3 Interpretability 163

164 Interpretability aligns model computations with 144
145 human cognition by providing explanations that are 146
147 not merely plausible but faithfully diagnostic (Guan 147
148 et al., 2025; Qin et al., 2025). Attribution analysis 149
150 quantifies token-level influence via $I(t_i) = \frac{\partial y}{\partial t_i}$ to 151
152 reveal causal contributions (Lanham et al., 2023; 153
154 Dziri et al., 2023). Complementarily, simulatability 155
156 evaluates whether humans can approximate $f_\theta(x)$ 157
158 from the reasoning trace R alone (Ye et al., 2023b; 159
160 Golovneva et al., 2023). High interpretability thus 161
162 enables epistemic alignment, rendering the model’s 163
164 latent reasoning cognitively reconstructable. 165

166 Formally, let $f_\theta(x) \rightarrow y$ denote the model 167
168 mapping. Interpretability requires a human- 169
170 comprehensible mapping $\Psi : R \mapsto E$, such that 171

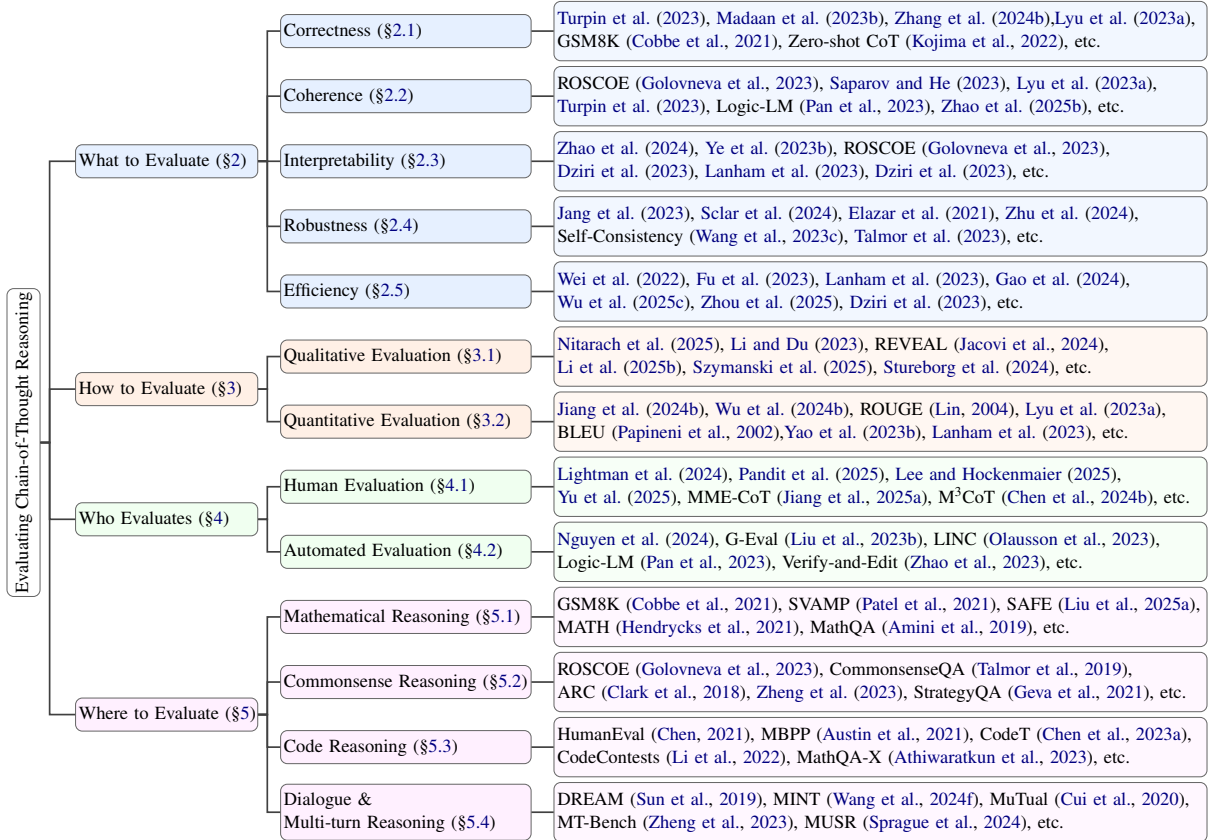


Figure 2: A brief taxonomy of representative works covered in our survey. The complete taxonomy is in Figure 6.

E faithfully captures the causal dependencies in f_θ . In other words, this formalizes the transparency and causal coherence of model reasoning.

2.4 Robustness

Robustness quantifies reasoning stability and reliability from random convergence by assessing the persistence of logical relations under paraphrasing or noise (Jang et al., 2023; Sclar et al., 2024). Typically, self-consistency metrics measure consensus across sampled reasoning paths (Wang et al., 2023c), while perturbation tests evaluate stability under adversarial or linguistic variations (Elazar et al., 2021; Zhu et al., 2024; Min et al., 2022).

For a model f_θ and perturbed input $x' \sim P(x)$, reasoning robustness is getting the most similar answer distributions:

$$\mathcal{R}(x) = 1 - \text{sim}(p(a|x), p(a|x')), \quad (2)$$

where sim is the similarity function of two probability distributions.

2.5 Efficiency

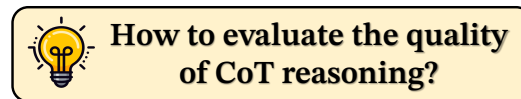
Efficiency reflects the optimal use of reasoning effort, maximizing progress per unit of computation. Practical assessment employs ratios of reasoning

length to solution depth, penalizing unnecessary verbosity (Lanham et al., 2023; Wu et al., 2025c; Wang and Zhou, 2024). Efficient reasoning thus achieves maximal insight with minimal resources.

Given a reasoning sequence $R = \{r_1, \dots, r_n\}$ with token cost $\tau(R)$ and performance or entropy $\mathcal{S}(R)$, the reasoning efficiency is quantified as

$$\mathcal{E}(R) = \frac{\mathcal{S}(R)}{\tau(R)}. \quad (3)$$

3 How to Evaluate



3.1 Qualitative Evaluation

3.1.1 Human Annotation

Human annotation provides gold-standard judgments for evaluating reasoning quality (Nitarach et al., 2025; Li and Du, 2023). To assess outputs systematically, define a fine-grained error taxonomy (e.g., factual errors, missing steps, and ambiguity) and verify each reasoning step to locate where CoT deviates from expert reasoning. Annotators also provide pairwise preferences (Zhang et al.,

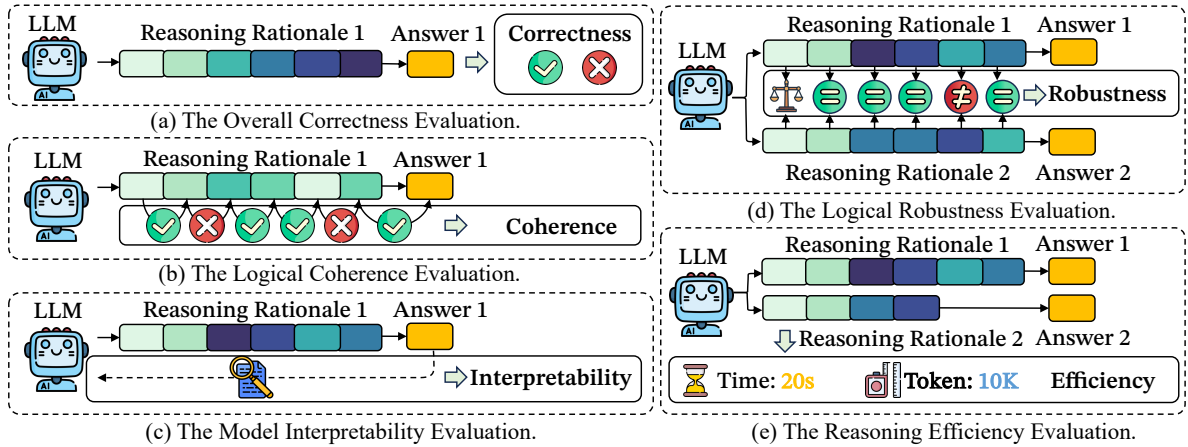


Figure 3: Five key dimensions of CoT evaluation: Correctness, Coherence, Interpretability, Robustness, and Efficiency. Each dimension captures a distinct facet of reasoning quality.

2024c; Zhao et al., 2025a; Yue et al., 2025), quantify agreement using Cohen’s κ or Krippendorff’s α and conduct periodic calibration to maintain consistency (Lee et al., 2024; Lyu et al., 2023b). Although costly to scale, this protocol yields reliable supervision for calibrating automated metrics and supporting robust model comparison (Chang et al., 2024; Chen et al., 2025e).

3.1.2 LLM-as-Judge

LLM-as-Judge uses advanced models to automate scalable reasoning assessment to provide a quality comment or select the superior reasoning trace between candidates, reducing calibration errors (Li et al., 2025b; Szymanski et al., 2025). For greater reliability, some frameworks augment qualitative critiques by CoT that anchor evaluations in explicit reasoning (Liu et al., 2023a; Jacovi et al., 2024). Step-level verification further breaks reasoning chains into atomic steps, checking each against formal constraints or external knowledge to detect hallucinations (Xu et al., 2024b; Paul et al., 2024). Despite the need for bias safeguards, LLM-as-Judge provides a scalable alternative to human experts for detailed reasoning evaluation.

3.2 Quantitative Evaluation

3.2.1 Correctness Evaluation

Quantitative metrics evaluate correctness by matching model outputs to gold-standard labels. For structured extraction, F1 and Exact Match test whether target tokens or error identifiers appear in the trace, enabling applications from pedagogical feedback (Jiang et al., 2024b) to debiasing interventions (Wu et al., 2024b). Recent studies use BLEU (Papineni et al., 2002) and ROUGE (Lin,

2004) to measure overlap between traces and reference explanations as a proxy for semantic similarity (Yao et al., 2023b; Maklad et al., 2025).

For subjective tasks, Likert ratings (1–5) score justification quality and are often combined with pairwise preferences against human reasoning (Bao et al., 2025; Jacovi et al., 2024; Li et al., 2025a). Discriminator-guided frameworks (Khalifa et al., 2023) and offline alignment systems (Wu et al., 2024a) train judge models to produce scores via ranking. Cog-CoT (Chen et al., 2025f) further uses hidden-state signals to predict trace reliability.

3.2.2 Coherence Evaluation

Intrinsic Evaluation assesses reasoning quality by separating the reasoning process from the final answer and focusing on structural and logical validity. It inspects intermediate steps for logical breaks, prioritizing derivation soundness over answer correctness (Golovneva et al., 2023; Lightman et al., 2024). Coherence is measured via step-wise analysis (Lyu et al., 2023a; Lanham et al., 2023), testing whether the reasoning chain causally supports the conclusion and separating deduction from heuristics (Saparov and He, 2023; Turpin et al., 2023). Building on this, structural analysis evaluates syntactic properties (e.g., step completeness and sequence length) to quantify inference validity and stability (Jiang et al., 2025b; Jin et al., 2024; Potamitis et al., 2025).

Recent methods capture reasoning integrity through logic-focused verification. VeriCoT (Feng et al., 2025) and DeduCE (Pandey et al., 2025) map language to symbolic logic for proof, while self-correction (Zhang et al., 2025a) and temporal checks (Mao et al., 2025) ensure consistency.

270	Causal metrics (Paul et al., 2024) and counterfactual tests (Wang et al., 2025b) distinguish inference from hallucinations. Granular tools include code unit tests (Saad-Falcon et al., 2024), logic probes (Kim et al., 2024a), and audits of multi-agent dynamics (Abdaljalil et al., 2025; She et al., 2025), shifting evaluation to structural audits.	321
271		322
272		323
273		324
274		325
275		326
276		327
277	Extrinsic Evaluation tests reasoning quality against external standards such as ground truth, structured knowledge, or execution environments (Lightman et al., 2024; Wang et al., 2024c; Liu et al., 2025a; Wang and Xu, 2025). This paradigm thus emphasizes objective verifiability rather than mere plausibility.	328
278		329
279		330
280		331
281		332
282		333
283		334
284		335
285	The most direct methods convert CoT into symbolic representations (Nguyen et al., 2025; Hu et al., 2025a,b) or executable code (Chen et al., 2022; Das et al., 2024), enabling deterministic checking (Olausson et al., 2023; Xu et al., 2024a; Feng et al., 2025; Simonds et al., 2025). LoT (Liu et al., 2025b) further integrate symbolic constraints dynamically into the reasoning process (Pan et al., 2023). Further, external compilers are applied to ensure language-agnostic verification (Nezhad et al., 2025). For knowledge-related reasoning, Graph-based methods (Luo et al., 2024; Amayuelas et al., 2025) trace reasoning paths through knowledge graphs, while fact decomposition (Min et al., 2023; Lage and Ostermann, 2025) and iterative feedback loops (Wang et al., 2024a) verify individual factual claims. For general generative tasks, semantic metrics (Sellam et al., 2020; Rei et al., 2020), like BERTScore (Zhang et al., 2020), quantify alignment with reference outputs, whereas structural metrics evaluate logical coherence (Prasad et al., 2023; Wang et al., 2025a; Hwang et al., 2025).	336
286		337
287		338
288		339
289		340
290		341
291		342
292		343
293		344
294		345
295		346
296		347
297		348
298		349
299		350
300		351
301		352
302		353
303		354
304		355
305		356
306		357
307		358
308		359
309		360
310		361
311		362
312		363
313		364
314		365
315		366
316		367
317		368
318		369
319		370
320		

via density metrics (Sui et al., 2025; Chen et al., 2025a). This entails calculating verbosity penalties and thought compression ratios (Saha et al., 2024; Ling et al., 2023b) to assess information density per token. Cost-effectiveness is evaluated by plotting accuracy-per-cost trade-offs, benchmarking frameworks like FrugalGPT (Chen et al., 2023b) and EcoAssistant (Zhang et al., 2023a) on dynamic query routing that invokes costly CoT paths only for complex cases, maximizing resource utility.

Takeaways

Prioritize Process: We should emphasize intrinsic validity and logical faithfulness to avoid “disguised accuracy”, where correct answers conceal flawed reasoning.

Hybrid Evaluation: We should integrate human/LLM judges for semantic nuance with external solvers for rigorous symbolic verification.

Resilient Efficiency: We should maintain stability against adversarial noise while balancing reasoning depth and efficiency.

4 Who Evaluates



Who evaluates the quality of the CoT?

4.1 Human Evaluation

4.1.1 Expert Evaluation

Expert evaluation provides ground truth. For process supervision, Lightman et al. (2024) employ annotators to flag step-level errors. Some domains require rigorous verification: Pandit et al. (2025) recruit mathematicians, while legal experts assess complex reasoning (Lee and Hockenmaier, 2025; Yu et al., 2025). In multimodal settings (Jiang et al., 2025a; Chen et al., 2024b), humans check visual-textual alignment. However, Tutek et al. (2025) caution that such approval reflects plausibility rather than parametric faithfulness.

4.1.2 Crowdsourced Evaluation

Crowdsourced evaluation engages broad populations to align reasoning metrics with general intuition. Specifically, Golovneva et al. (2023) have workers validate error typologies, matching metrics like semantic consistency to layperson judgments. Likewise, Ramnath et al. (2024) gather large-scale feedback for MaRio’s optimization, confirming

preferences for plausible rationales. However, non-experts reduce precision. Kumar et al. (2025) show crowd workers miss subtle nuances versus domain experts, especially in complex tasks.

4.1.3 Human-In-The-Loop Evaluation

Human-In-The-Loop evaluation balances scalability and precision by directing human effort to key verification steps. Specifically, CRV (Zhao et al., 2025b) has humans validate model-flagged error fingerprints, avoiding full output review. To optimize effort, Diao et al. (2024) focus annotators on high-uncertainty samples in Active-Prompt, using human judgment to resolve ambiguity. Likewise, Li et al. (2023a) place humans as the final quality gate in HaluEval, with annotators conducting sanity checks only after programmatic filtering to minimize manual work.

4.2 Automated Evaluation

4.2.1 Symbolic and Rule-Based Verification

Symbolic and rule-based verification provides deterministic checks. For factuality, knowledge graphs anchor CoT to structured entities to evaluate hallucinations (Nguyen et al., 2024; Zhao et al., 2023; Peng et al., 2023). For logical rigor, solvers like LINC (Olausson et al., 2023), Logic-LM (Pan et al., 2023), and VeriCoT (Feng et al., 2025) convert steps to formal proofs, while theorem provers constrain generation to deductive rules (Poesia et al., 2023; Ye et al., 2023a; Mysore et al., 2023; Ling et al., 2023a).

4.2.2 Encoder-Based Verification

Encoder-based verification evaluates rationales via distance metrics against golden chains. While semantic encoders like BERTScore (Zhang et al., 2020) and BLEURT (Sellam et al., 2020) measure similarity, they often overlook reasoning rigor. Addressing this, ROSCOE (Golovneva et al., 2023) assesses step-wise consistency, whereas RecEval (Prasad et al., 2023) employs natural language inference to verify strict logical entailment and mitigate shortcut learning.

4.2.3 Reward-Based Verification

Reward-based verification uses learned models to score reasoning quality (Lightman et al., 2024; Duan et al., 2025; Wang et al., 2025c). Process Reward Models (PRMs) act as step-by-step supervisors (Lightman et al., 2024; Wang et al., 2024c; Song et al., 2025), detecting intermediate errors overlooked by outcome metrics. Fine-grained detectors in REVEAL (Jacovi et al., 2024) and Li

et al. (2024b) categorize failures (e.g., calculation errors, symbol misuse), mimicking expert diagnostics. These PRMs evaluate logical transitions at each step (Uesato et al., 2022), reducing disguised accuracy (Bentham et al., 2024) and ensuring sound reasoning trajectories.

4.2.4 Generative Model Verification

Generative verification employs LLMs as judges to critique reasoning via semantic understanding. G-Eval (Liu et al., 2023b) uses GPT-4 for reference-free evaluation. Prompt-based systems like Self-Refine (Madaan et al., 2023b) and Self-Reflection (Shinn et al., 2023; Ji et al., 2023b; Weng et al., 2023) enable iterative self-correction of hallucinations. Recent benchmarks (Li et al., 2025a,b; Zhang et al., 2024c; Lee et al., 2025b; Chen et al., 2025c) assess these capabilities for expert-level rigor, while Zhang et al. (2025b) integrates verification into generation (Ali et al., 2025; Xu et al., 2025) to improve inference chains.

Takeaways

Automation Does Not Guarantee Objectivity: Automated methods are efficient but carry model biases and shortcuts.

Hybrid Intelligence as Next Paradigm: Human-AI evaluation is a distributed cognitive architecture. Future systems will exploit these complementarities.

Verification as Iterative Feedback: Evaluation evolves from post-hoc measurement to a core learning signal that refines alignment and reasoning fidelity.

5 Where to Evaluate



Where is the CoT quality evaluation performed?

We review key evaluation aspects across common CoT scenarios here. Representative datasets and evaluation methodologies are in Appendix D.

5.1 Mathematical Reasoning

Mathematical reasoning failures range from calculation errors within coherent logic (e.g., Minerva (Lewkowycz et al., 2022)) to reasoning failures, such as invalid deductions or knowledge gaps, that undermine validity (Jiang et al., 2024b; Golovneva et al., 2023). Crucially, outcome accuracy often diverges from process validity, masking poor reasoning behind high scores (Guo et al.,

2025b). While PRMs address this via consistency checks, their reliability remains limited by noisy reward estimation (Zhang et al., 2025d).

5.2 Commonsense Reasoning

Commonsense reasoning evaluation typically uses external verification against knowledge bases like ConceptNet (Liu and Singh, 2004) and ATOMIC (Sap et al., 2019). Methods always align intermediate steps with explicit knowledge paths, ensuring Correctness via factual grounding rather than plausibility alone (Feng et al., 2024; Wang et al., 2024a). However, quantifying logical coherence is challenging and still heavily relies on human judgment, as commonsense is probabilistic and allows multiple valid traces (Miao et al., 2024).

5.3 Code Reasoning

Code reasoning functions as a deterministic sandbox, verifying validity via execution (Chen, 2021). Evaluation should extend beyond correctness to capture efficiency and robustness. Benchmarks like APPS (Hendrycks et al., 2021) incorporate resource constraints, while CodeT (Chen et al., 2023a) employs dual-execution to reveal latent failures overlooked by standard checks.

5.4 Dialogue and Multi-turn Reasoning

In dialogue and multi-turn settings, logical coherence hinges on interpretability and reliable coreference resolution amid implicit assumptions. Mutual (Cui et al., 2020) stresses history-based inference, and QReCC (Anantha et al., 2021) assesses state tracking through context-dependent query rewriting. Coherence must persist across interactions: DNLI (Welleck et al., 2019) classifies response relations as entailment, neutral, or contradiction, while CICERO (Ghosal et al., 2022) penalizes hallucinations that disrupt causal continuity and conversational flow.²

Takeaways

Different Evaluation Scopes: Evaluation scenarios vary in CoT reasoning demands.

External Grounding: Validate reasoning with executable suites or knowledge graphs, beyond surface plausibility.

Contextual Consistency: In long-context scenarios, ensure reliability by tracking implicit assumptions and maintaining causal coherence across turns.

²Scientific reasoning is discussed in Appendix D.5.

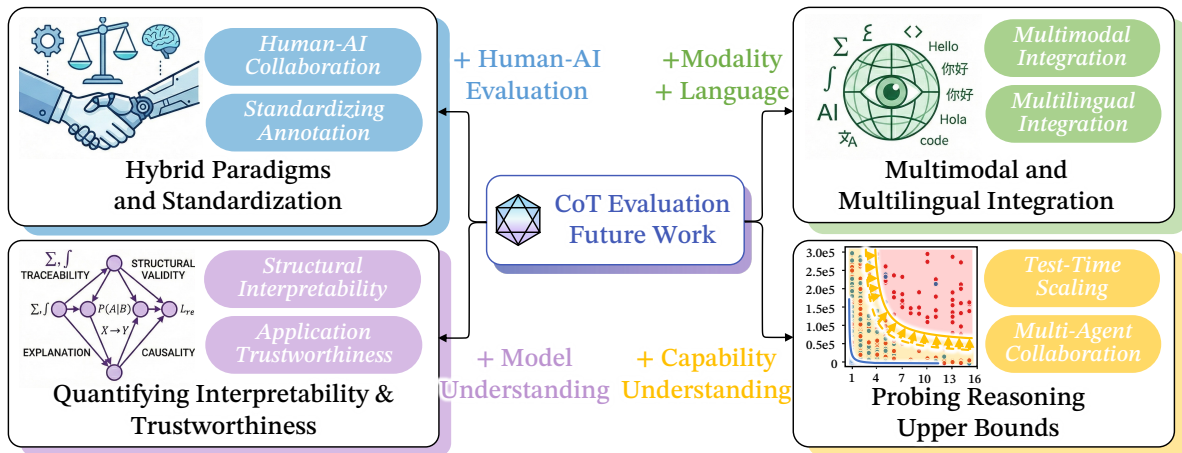


Figure 4: Future Directions for CoT Evaluation. We identify four key directions to advance the field.

6 Future Directions

As LLMs evolve from chat agents into specialized reasoning engines, CoT evaluation should transform in parallel to ensure rigorous assessment. We identify four critical frontiers (illustrated in Figure 4) essential to bridge current methodological gaps. For a comprehensive discussion of these strategic research trajectories, please refer to Appendix E.

6.1 Hybrid Paradigms and Standardization

Future protocols should prioritize **human-AI collaboration**, using models for initial filtering and humans for targeted verification (Diao et al., 2024; Li et al., 2023a). **Standardizing annotation** guidelines and testing adversarial robustness are crucial for scalable assessment (Cui et al., 2025). Thus, research must develop techniques to generate adversarial samples that stress-test metrics against subtle prompt variations.

6.2 Multimodal and Multilingual Frontiers

Evaluation should extend beyond text-only settings to multimodal models (Jiang et al., 2025a; Wu et al., 2025a) and **multilingual** CoT reasoning (Qin et al., 2023), with explicit tests of cross-modal grounding and cross-lingual logical consistency. To reduce Anglocentric bias and to better characterize **multimodal** capabilities, benchmarks should adopt more fine-grained protocols that cover diverse languages and require deeper visual-text reasoning (Bang et al., 2023; Chen et al., 2024b).

6.3 Quantifying Interpretability & Trustworthiness

Research should quantify **structural interpretability** by linking CoT traces to causal or knowledge graphs (Amayuelas et al., 2025; Nguyen et al.,

2024) and using XAI to attribute input-context influence per reasoning step (Turpin et al., 2023). In law and healthcare, evaluations should prioritize **application trustworthiness** over accuracy (Yu et al., 2025). Enterprise frameworks should therefore enable interpretability, bias detection, and auditability (Shaikh et al., 2023).

6.4 Probing Reasoning Upper Bounds

Current benchmarks largely measure static inference, overlooking the potential of **test-time scaling** where models simulate deliberate thinking. In this paradigm, extra compute translates into verified gains through extended reasoning steps or deeper search (Chen et al., 2024a, 2025b; Snell et al., 2024; Muennighoff et al., 2025; Zhang et al., 2025c; Chen et al., 2025d). Similarly, **multi-agent collaboration** leverages ensemble diversity to estimate swarm-intelligence upper bounds (Du et al., 2023; Guo et al., 2024; Qian et al., 2024). Future evaluation should quantify the efficiency-performance trade-off, measuring how reasoning capabilities evolve as computing and collaboration scale.

7 Conclusion

This survey demonstrates that evaluating reasoning differs fundamentally from scoring outcomes. Transcending accuracy as a weak proxy, we structure assessment along five core dimensions to frame evaluation as a multifaceted design space. Reliable verification demands a synergy of symbolic rigor, scalable automation, and human judgment rather than a single paradigm. Future progress relies on multimodal, multilingual hybrids for trustworthy reasoning bounds. We hope this work offers insightful perspectives and foundational guidance for advancing the CoT reasoning evaluation field.

596 Limitations

597 Despite our efforts to provide a systematic and
598 comprehensive survey, this work faces several limi-
599 tations inherent to the rapid evolution of the field.

600 **Scope Delimitation.** To maintain depth, we
601 strictly narrowed our scope to the *evaluation*
602 of Chain-of-Thought reasoning. While train-
603 ing methodologies (e.g., Process Reward Models,
604 RLHF) inherently involve evaluation steps, we ex-
605 cluded papers that focus solely on optimizing train-
606 ing objectives without introducing distinct assess-
607 ment metrics or frameworks. This boundary, while
608 necessary for focus, excludes broader alignment
609 techniques that indirectly impact reasoning fidelity.

610 **Taxonomical Fluidity.** The proposed “What-
611 How-Who-Where” framework serves as a struc-
612 tured lens to organize fragmented literature. In
613 practice, however, these dimensions are not mutu-
614 ally exclusive. Many hybrid systems span multiple
615 categories (e.g., a generative judge employing rule-
616 based constraints), and our categorization reflects
617 their primary methodological contribution rather
618 than a rigid classification.

619 References

620 Samir Abdaljalil, Hasan Kurban, Khalid Qaraqe, and
621 Erchin Serpedin. 2025. [Theorem-of-thought: A
622 multi-agent framework for abductive, deductive, and
623 inductive reasoning in language models.](#) *ArXiv
624 preprint*, abs/2506.07106.

625 Nassir Jabir Al-Khafaji and Basit Khalaf Majeed.
626 2024. Evaluating large language models using arabic
627 prompts to generate python codes. In *2024 4th Inter-
628 national Conference on Emerging Smart Technolo-
629 gies and Applications (eSmarTA)*, pages 1–5. IEEE.

630 Yusuf Ali, Gryphon Patlin, Karthik Kothuri, Muham-
631 mad Zubair Irshad, Wuwei Liang, and Zsolt Kira.
632 2025. Eve: A generator-verifier system for genera-
633 tive policies. *arXiv preprint arXiv:2512.21430*.

634 Aylton Almeida, Laerte Xavier, and Marco Tulio Va-
635 lente. 2024. Automatic library migration using large
636 language models: First results. In *Proceedings of the
637 18th ACM/IEEE International Symposium on Empiri-
638 cal Software Engineering and Measurement*, pages
639 427–433.

640 Alfonso Amayuelas, Joy Sain, Simerjot Kaur, and
641 Charese Smiley. 2025. [Grounding llm reasoning with
642 knowledge graphs.](#) *ArXiv preprint*, abs/2502.13247.

643 Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik
644 Koncel-Kedziorski, Yejin Choi, and Hannaneh Ha-
645 jishirzi. 2019. [MathQA: Towards interpretable math](#)

[word problem solving with operation-based for-
646 malisms.](#) In *Proceedings of the 2019 Conference
647 of the North American Chapter of the Association for
648 Computational Linguistics: Human Language Tech-
649 nologies, Volume 1 (Long and Short Papers)*, pages
650 2357–2367, Minneapolis, Minnesota. Association for
651 Computational Linguistics. 652

Raviteja Anantha, Svitlana Vakulenko, Zhucheng Tu,
Shayne Longpre, Stephen Pulman, and Srinivas
Chappidi. 2021. [Open-domain question answering
653 goes conversational via question rewriting.](#) In *Pro-
654 ceedings of the 2021 Conference of the North Amer-
655 ican Chapter of the Association for Computational
656 Linguistics: Human Language Technologies*, pages
657 520–534, Online. Association for Computational Lin-
658 guistics. 659

Pepa Atanasova, Oana-Maria Camburu, Christina Li-
oma, Thomas Lukasiewicz, Jakob Grue Simonsen,
and Isabelle Augenstein. 2023. [Faithfulness tests
660 for natural language explanations.](#) In *Proceedings
661 of the 61st Annual Meeting of the Association for
662 Computational Linguistics (Volume 2: Short Papers)*,
663 pages 283–294, Toronto, Canada. Association for
664 Computational Linguistics. 665

Ben Athiwaratkun, Sanjay Krishna Gouda, Zijian Wang,
Xiaopeng Li, Yuchen Tian, Ming Tan, Wasi Uddin
Ahmad, Shiqi Wang, Qing Sun, Mingyue Shang, Su-
jan Kumar Gonugondla, Hantian Ding, Varun Ku-
mar, Nathan Fulton, Arash Farahani, Siddhartha Jain,
Robert Giaquinto, Haifeng Qian, Murali Krishna
Ramanathan, and Ramesh Nallapati. 2023. [Multi-
666 lingual evaluation of code generation models.](#) In *The
667 Eleventh International Conference on Learning Rep-
668 resentations, ICLR 2023, Kigali, Rwanda, May 1-5,
669 2023*. OpenReview.net. 670

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten
Bosma, Henryk Michalewski, David Dohan, Ellen
Jiang, Carrie Cai, Michael Terry, Quoc Le, and 1
others. 2021. [Program synthesis with large language
671 models.](#) *ArXiv preprint*, abs/2108.07732. 672

Sriram Balasubramanian, Samyadeep Basu, and So-
heil Feizi. 2025. [A closer look at bias and chain-of-
673 thought faithfulness of large \(vision\) language mod-
674 els.](#) *ArXiv preprint*, abs/2505.23945. 675

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-
liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei
Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu,
and Pascale Fung. 2023. [A multitask, multilingual,
676 multimodal evaluation of ChatGPT on reasoning, hal-
677 lucination, and interactivity.](#) In *Proceedings of the
678 13th International Joint Conference on Natural Lan-
679 guage Processing and the 3rd Conference of the Asia-
680 Pacific Chapter of the Association for Computational
681 Linguistics (Volume 1: Long Papers)*, pages 675–718,
682 Nusa Dua, Bali. Association for Computational Lin-
683 guistics. 684

Guangsheng Bao, Hongbo Zhang, Cunxiang Wang,
Linyi Yang, and Yue Zhang. 2025. How likely do
685

704	llms with cot mimic human reasoning? In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 7831–7850.	A survey on evaluation of large language models. <i>ACM transactions on intelligent systems and technology</i> , 15(3):1–45.	761
705			762
706			763
707	Oliver Bentham, Nathan Stringham, and Ana Marasović. 2024. Chain-of-thought unfaithfulness as disguised accuracy . <i>ArXiv preprint</i> , abs/2402.14897.	Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2023a. Codet: Code generation with generated tests . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	764
708			765
709			766
710	Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. 2024. Graph of thoughts: Solving elaborate problems with large language models . In <i>Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada</i> , pages 17682–17690. AAAI Press.	Lingjiao Chen, Matei Zaharia, and James Zou. 2023b. Frugalgpt: How to use large language models while reducing cost and improving performance . <i>ArXiv preprint</i> , abs/2305.05176.	767
711			768
712			769
713			770
714			771
715			772
716			773
717			774
718			775
719			776
720			777
721			778
722	Ning Bian, Xianpei Han, Le Sun, Hongyu Lin, Yaojie Lu, Ben He, Shanshan Jiang, and Bin Dong. 2024. ChatGPT is a knowledgeable but inexperienced solver: An investigation of commonsense problem in large language models . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 3098–3110, Torino, Italia. ELRA and ICCL.	Qiguang Chen, Dengyun Peng, Jinhao Liu, HuiKang Su, Jiannan Guan, Libo Qin, and Wanxiang Che. 2025a. Aware first, think less: Dynamic boundary self-awareness drives extreme reasoning efficiency in large language models . <i>arXiv preprint arXiv:2508.11582</i> .	779
723			780
724			781
725			782
726			783
727			784
728			785
729			786
730			787
731	Rishi Bommasani and 1 others. 2021. On the opportunities and risks of foundation models . <i>ArXiv preprint</i> , abs/2108.07258.	Qiguang Chen, Libo Qin, Jinhao Liu, Yue Liao, Jiaqi Wang, Jingxuan Zhou, and Wanxiang Che. 2025b. Rbf++: Quantifying and optimizing reasoning boundaries across measurable and unmeasurable capabilities for chain-of-thought reasoning . <i>arXiv preprint arXiv:2505.13307</i> .	788
732			789
733			790
734	Ana Brassard, Benjamin Heinzerling, Keito Kudo, Keisuke Sakaguchi, and Kentaro Inui. 2024. Acorn: Aspect-wise commonsense reasoning explanation evaluation . <i>ArXiv preprint</i> , abs/2405.04818.	Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. 2025c. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models . <i>ArXiv preprint</i> , abs/2503.09567.	791
735			792
736			793
737			794
738	Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. 2024. Large language monkeys: Scaling inference compute with repeated sampling . <i>arXiv preprint arXiv:2407.21787</i> .	Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiaqi Wang, Mengkang Hu, Zhi Chen, Wanxiang Che, and Ting Liu. 2025d. Ecm: A unified electronic circuit model for explaining the emergence of in-context learning and chain-of-thought in large language model . <i>arXiv preprint arXiv:2502.03325</i> .	795
739			796
740			797
741			798
742			799
743	Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, and 1 others. 2023. Open problems and fundamental limitations of reinforcement learning from human feedback . <i>ArXiv preprint</i> , abs/2307.15217.	Qiguang Chen, Libo Qin, Jiaqi Wang, Jingxuan Zhou, and Wanxiang Che. 2024a. Unlocking the capabilities of thought: A reasoning boundary framework to quantify and optimize chain-of-thought . <i>Advances in Neural Information Processing Systems</i> , 37:54872–54904.	800
744			801
745			802
746			803
747			804
748			805
749			806
750	Kranti Chalamalasetti, Jana Götze, Sherzod Hakimov, Brielen Madureira, Philipp Sadler, and David Schlangen. 2023. clembench: Using game play to evaluate chat-optimized language models as conversational agents . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 11174–11219, Singapore. Association for Computational Linguistics.	Qiguang Chen, Libo Qin, Jin Zhang, Zhi Chen, Xiao Xu, and Wanxiang Che. 2024b. M³cot: A novel benchmark for multi-domain multi-step multi-modal chain-of-thought . In <i>Proc. of ACL</i> .	807
751			808
752			809
753			810
754			811
755			812
756			813
757			814
758	Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, and 1 others. 2024.	Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. 2022. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks . <i>ArXiv preprint</i> , abs/2211.12588.	815
759			816
760			817

815	Xiaoshu Chen, Sihang Zhou, Ke Liang, Duanyang Yuan,	Mingkai Deng, Bowen Tan, Zhengzhong Liu, Eric Xing,	870
816	Haoyuan Chen, Xiaoyu Sun, Linyuan Meng, and	and Zhiting Hu. 2021. Compression, transduction,	871
817	Xinwang Liu. 2025e. Putting on the thinking hats:	and creation: A unified framework for evaluating	872
818	A survey on chain of thought fine-tuning from the	natural language generation. In <i>Proceedings of the</i>	873
819	perspective of human reasoning mechanism. <i>arXiv</i>	<i>2021 Conference on Empirical Methods in Natural</i>	874
820	preprint arXiv:2510.13170.	<i>Language Processing</i> , pages 7580–7605, Online and	875
		Punta Cana, Dominican Republic. Association for	876
821	Xinyun Chen, Renat Aksitov, Uri Alon, Jie Ren, Ke-	Computational Linguistics.	877
822	fan Xiao, Pengcheng Yin, Sushant Prakash, Charles		
823	Sutton, Xuezhi Wang, and Denny Zhou. 2023c. Uni-	Xun Deng, Sicheng Zhong, Honghua Dong, Jingyu	878
824	versal self-consistency for large language model	Hu, Sidi Mohamed Beillahi, Xujie Si, and Fan Long.	879
825	generation. <i>ArXiv preprint</i> , abs/2311.17311.	2024. Assessing code generation with intermediate	880
		languages. <i>ArXiv preprint</i> , abs/2407.05411.	881
826	Zhipeng Chen, Kun Zhou, Beichen Zhang, Zheng Gong,	Shizhe Diao, Pengcheng Wang, Yong Lin, Rui Pan, Xi-	882
827	Xin Zhao, and Ji-Rong Wen. 2023d. ChatCoT:	ang Liu, and Tong Zhang. 2024. Active prompting	883
828	Tool-augmented chain-of-thought reasoning on chat-	with chain-of-thought for large language models. In	884
829	based large language models. In <i>Findings of the</i>	<i>Proceedings of the 62nd Annual Meeting of the As-</i>	885
830	<i>Association for Computational Linguistics: EMNLP</i>	<i>sociation for Computational Linguistics (Volume 1:</i>	886
831	2023, pages 14777–14790, Singapore. Association	<i>Long Papers)</i> , pages 1330–1350.	887
832	for Computational Linguistics.		
833	Zijun Chen, Wenbo Hu, and Richang Hong. 2025f.	Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenen-	888
834	Deep hidden cognition facilitates reliable chain-of-	baum, and Igor Mordatch. 2023. Improving factual-	889
835	thought reasoning. <i>ArXiv preprint</i> , abs/2507.10007.	ity and reasoning in language models through multi-	890
		agent debate. In <i>Forty-first International Conference</i>	891
836	Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang	<i>on Machine Learning.</i>	892
837	Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu,	Keyu Duan, Zichen Liu, Xin Mao, Tianyu Pang,	893
838	Bing Qin, and Ting Liu. 2024. Navigate through enig-	Changyu Chen, Qiguang Chen, Michael Qizhe Shieh,	894
839	matic labyrinth a survey of chain of thought reason-	and Longxu Dou. 2025. Efficient process reward	895
840	ing: Advances, frontiers and future. In <i>Proceedings</i>	model training via active learning. <i>arXiv preprint</i>	896
841	<i>of the 62nd Annual Meeting of the Association for</i>	<i>arXiv:2504.10559.</i>	897
842	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	Zackary Okun Dunivin. 2024. Scalable qualitative cod-	898
843	pages 1173–1203.	ing with llms: Chain-of-thought reasoning matches	899
844	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,	human performance in some hermeneutic tasks.	900
845	Ashish Sabharwal, Carissa Schoenick, and Oyvind	<i>ArXiv preprint</i> , abs/2401.15170.	901
846	Tafjord. 2018. Think you have solved question		
847	answering? try arc, the ai2 reasoning challenge. <i>ArXiv</i>	Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine	902
848	preprint , abs/1803.05457.	Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, Peter	903
849	Karl Cobbe and 1 others. 2021. Training verifiers to	West, Chandra Bhagavatula, Ronan Le Bras, Jena D.	904
850	solve math word problems. volume abs/2110.14168.	Hwang, Soumya Sanyal, Xiang Ren, Allyson Et-	905
851	Leyang Cui, Yu Wu, Shujie Liu, Yue Zhang, and Ming	tinger, Zaid Harchaoui, and Yejin Choi. 2023. Faith	906
852	Zhou. 2020. MuTual: A dataset for multi-turn dia-	and fate: Limits of transformers on composition-	907
853	logue reasoning. In <i>Proceedings of the 58th Annual</i>	ality. In <i>Advances in Neural Information Processing</i>	908
854	<i>Meeting of the Association for Computational Lin-</i>	<i>Systems 36: Annual Conference on Neural Informa-</i>	909
855	<i>guistics</i> , pages 1406–1416, Online. Association for	<i>tion Processing Systems 2023, NeurIPS 2023, New</i>	910
856	Computational Linguistics.	<i>Orleans, LA, USA, December 10 - 16, 2023.</i>	911
857	Yu Cui, Bryan Hooi, Yujun Cai, and Yiwei Wang. 2025.	Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhi-	912
858	Process or result? manipulated ending tokens can	lasha Ravichander, Eduard Hovy, Hinrich Schütze,	913
859	mislead reasoning llms to ignore the correct reason-	and Yoav Goldberg. 2021. Measuring and improving	914
860	ing steps. <i>ArXiv preprint</i> , abs/2503.19326.	consistency in pretrained language models. <i>Transac-</i>	915
		<i>tions of the Association for Computational Linguis-</i>	916
861	Debrup Das, Debopriyo Banerjee, Somak Aditya, and	<i>tics</i> , 9:1012–1031.	917
862	Ashish Kulkarni. 2024. MATHSENSEI: A tool-	Run-Ze Fan, Zengzhi Wang, and Pengfei Liu. 2025.	918
863	augmented large language model for mathematical	Megascience: Pushing the frontiers of post-training	919
864	reasoning. In <i>Proceedings of the 2024 Conference of</i>	datasets for science reasoning. <i>ArXiv preprint</i> ,	920
865	<i>the North American Chapter of the Association for</i>	abs/2507.16812.	921
866	<i>Computational Linguistics: Human Language Tech-</i>	Hao Fei, Bobo Li, Qian Liu, Lidong Bing, Fei Li, and	922
867	<i>nologies (Volume 1: Long Papers)</i> , pages 942–966,	Tat-Seng Chua. 2023. Reasoning implicit sentiment	923
868	Mexico City, Mexico. Association for Computational	with chain-of-thought prompting. In <i>Proceedings</i>	924
869	Linguistics.	<i>of the 61st Annual Meeting of the Association for</i>	925
		<i>Computational Linguistics (Volume 2: Short Papers)</i> ,	926

927	pages 1171–1182, Toronto, Canada. Association for Computational Linguistics.	
928		
929	Kehua Feng, Keyan Ding, Weijie Wang, Xiang Zhuang, Zeyuan Wang, Ming Qin, Yu Zhao, Jianhua Yao, Qiang Zhang, and Huajun Chen. 2024. Sciknoweval: Evaluating multi-level scientific knowledge of large language models . <i>ArXiv preprint</i> , abs/2406.09098.	
930		
931		
932		
933		
934	Yu Feng, Nathaniel Weir, Kaj Bostrom, Sam Bayless, Darion Cassel, Sapana Chaudhary, Benjamin Kiesl-Reiter, and Huzefa Rangwala. 2025. Vericot: Neuro-symbolic chain-of-thought validation via logical consistency checks . <i>ArXiv preprint</i> , abs/2511.04662.	
935		
936		
937		
938		
939	Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. 2023. Complexity-based prompting for multi-step reasoning . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
940		
941		
942		
943		
944	Peng Gao and 1 others. 2024. Ruler: What’s the real reasoning capability of your language model? <i>ArXiv preprint</i> , abs/2402.18679.	
945		
946		
947	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies . <i>Transactions of the Association for Computational Linguistics</i> , 9:346–361.	
948		
949		
950		
951		
952		
953	Deepanway Ghosal, Siqi Shen, Navonil Majumder, Rada Mihalcea, and Soujanya Poria. 2022. CICERO: A dataset for contextualized commonsense inference in dialogues . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5010–5028, Dublin, Ireland. Association for Computational Linguistics.	
954		
955		
956		
957		
958		
959		
960	Olga Golovneva, Moya Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. 2023. ROSCOE: A suite of metrics for scoring step-by-step reasoning . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
961		
962		
963		
964		
965		
966		
967	Google DeepMind. 2025. Gemini 3 models and ecosystem . Online.	
968		
969	Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. 2023. Think before you speak: Training language models with pause tokens . <i>arXiv preprint arXiv:2310.02226</i> .	
970		
971		
972		
973		
974	Alex Gu, Baptiste Rozière, Hugh James Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida Wang. 2024. Cruxeval: A benchmark for code reasoning, understanding and execution . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	
975		
976		
977		
978		
979		
980		
	Melody Y. Guan, Miles Wang, Micah Carroll, Zehao Dou, Annie Y. Wei, Marcus Williams, Benjamin Arnav, Joost Huizinga, Ian Kivlichan, Mia Glaese, Jakub Pachocki, and Bowen Baker. 2025. Monitoring monitorability .	981
		982
		983
		984
		985
	Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and Dawn Song. 2023. The false promise of imitating proprietary llms . <i>ArXiv preprint</i> , abs/2305.15717.	986
		987
		988
		989
	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025a. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning . <i>ArXiv preprint</i> , abs/2501.12948.	990
		991
		992
		993
		994
	Jiaying Guo, Wenjie Yang, Shengzhong Zhang, Tongshan Xu, Lun Du, Da Zheng, and Zengfeng Huang. 2025b. Right is not enough: The pitfalls of outcome supervision in training llms for math reasoning . <i>ArXiv preprint</i> , abs/2506.06877.	995
		996
		997
		998
		999
	Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xi-angliang Zhang. 2024. Large language model based multi-agents: A survey of progress and challenges . <i>arXiv preprint arXiv:2402.01680</i> .	1000
		1001
		1002
		1003
		1004
	Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Linhao Yu, Yan Liu, Jiaxuan Li, Bojian Xiong, Deyi Xiong, and 1 others. 2023. Evaluating large language models: A comprehensive survey . <i>ArXiv preprint</i> , abs/2310.19736.	1005
		1006
		1007
		1008
		1009
	Zhiyuan He and Dingmin Wang. 2025. When do symbolic solvers enhance reasoning in large language models? <i>ArXiv preprint</i> , abs/2512.03272.	1010
		1011
		1012
	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset . <i>ArXiv preprint</i> , abs/2103.03874.	1013
		1014
		1015
		1016
		1017
	Or Honovich, Leshem Choshen, Roei Aharoni, Ella Neeman, Idan Szpektor, and Omri Abend. 2021. q^2: Evaluating factual consistency in knowledge-grounded dialogues via question generation and question answering . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7856–7870, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	1018
		1019
		1020
		1021
		1022
		1023
		1024
		1025
		1026
	Md Sifat Hossain, Anika Tabassum, Md Fahim Arfin, and Tarannum Shaila Zaman. 2025. Llm-pros: Analyzing large language models’ performance in competitive problem solving . In <i>2025 IEEE/ACM International Workshop on Large Language Models for Code (LLM4Code)</i> , pages 80–87. IEEE.	1027
		1028
		1029
		1030
		1031
		1032
	Mengkang Hu, Tianxing Chen, Yude Zou, Yuheng Lei, Qiguang Chen, Ming Li, Yao Mu, Hongyuan Zhang, Wenqi Shao, and Ping Luo. 2025a. Text2world :	1033
		1034
		1035

1036	Benchmarking large language models for symbolic world model generation. <i>arXiv preprint arXiv:2502.13092</i> .	case study with negated prompts. In <i>Transfer learning for natural language processing workshop</i> , pages 52–62. PMLR.	1092
1037			1093
1038			1094
1039	Mengkang Hu, Bowei Xia, Yuran Wu, Ailing Yu, Yude Zou, Qiguang Chen, Shijian Wang, Jiarui Jin, Kexin Li, Wenxiang Jiao, and 1 others. 2025b. Agent2world: Learning to generate symbolic world models via adaptive multi-agent feedback. <i>arXiv preprint arXiv:2512.22336</i> .	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023a. Survey of hallucination in natural language generation. <i>ACM computing surveys</i> , 55(12):1–38.	1095
1040			1096
1041			1097
1042			1098
1043			1099
1044			
1045	Jie Huang and Kevin Chang. 2023. Large language models as general pattern machines . <i>ArXiv preprint</i> , abs/2306.09831.	Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. 2023b. Towards mitigating hallucination in large language models via self-reflection. <i>arXiv preprint arXiv:2310.06271</i> .	1100
1046			1101
1047			1102
1048			1103
1049	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and 1 others. 2025a. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. <i>ACM Transactions on Information Systems</i> , 43(2):1–55.	Dongzhi Jiang, Renrui Zhang, Ziyu Guo, Yanwei Li, Yu Qi, Xinyan Chen, Liuhui Wang, Jianhan Jin, Claire Guo, Shen Yan, and 1 others. 2025a. Mme-cot: Benchmarking chain-of-thought in large multimodal models for reasoning quality, robustness, and efficiency . <i>ArXiv preprint</i> , abs/2502.09621.	1104
1050			1105
1051			1106
1052			1107
1053			1108
1054			1109
1055	Zhen Huang, Zengzhi Wang, Shijie Xia, Xuefeng Li, Haoyang Zou, Ruijie Xu, Run-Ze Fan, Lyumanshan Ye, Ethan Chern, Yixin Ye, Yikai Zhang, Yuqing Yang, Ting Wu, Binjie Wang, Shichao Sun, Yang Xiao, Yiyuan Li, Fan Zhou, Steffi Chern, and 9 others. 2024. Olympicarena: Benchmarking multi-discipline cognitive reasoning for superintelligent AI . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	Gangwei Jiang, Yahui Liu, Zhaoyi Li, Wei Bi, Fuzheng Zhang, Linqi Song, Ying Wei, and Defu Lian. 2025b. What makes a good reasoning chain? uncovering structural patterns in long chain-of-thought reasoning. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 6501–6525.	1110
1056			1111
1057			1112
1058			1113
1059			1114
1060			1115
1061			1116
1062			
1063			1117
1064			1118
1065			1119
1066	Zhiyuan Huang, Baichuan Yang, Zikun He, Yanhong Wu, Fang Hongyu, Zhenhe Liu, Lin Dongsheng, and Bing Su. 2025b. Chemvts-bench: Evaluating visual-textual-symbolic reasoning of multimodal large language models in chemistry . <i>ArXiv preprint</i> , abs/2511.17909.	Yuan Jiang, Yujian Zhang, Liang Lu, Christoph Treude, Xiaohong Su, Shan Huang, and Tiantian Wang. 2025c. Enhancing high-quality code generation in large language models with comparative prefix-tuning . <i>ArXiv preprint</i> , abs/2503.09020.	1120
1067			1121
1068			
1069			1122
1070			1123
1071			1124
1072			1125
1073	Hyeon Hwang, Yewon Cho, Chanwoong Yoon, Yein Park, Minju Song, Kyungjae Lee, Gangwoo Kim, and Jaewoo Kang. 2025. Assessing llm reasoning steps via principal knowledge grounding. In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 19925–19948.	Zheng Ping Jiang, Yining Lu, Hanjie Chen, Daniel Khashabi, Benjamin Van Durme, and Anqi Liu. 2024a. Rora: Robust free-text rationale evaluation. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1070–1087.	1126
1074			1127
1075			
1076			1128
1077			1129
1078	Alon Jacovi, Yonatan Bitton, Bernd Bohnet, Jonathan Herzig, Or Honovich, Michael Tseng, Michael Collins, Roei Aharoni, and Mor Geva. 2024. A chain-of-thought is as strong as its weakest link: A benchmark for verifiers of reasoning chains . <i>ArXiv preprint</i> , abs/2402.00559.	Zhuoxuan Jiang, Haoyuan Peng, Shanshan Feng, Fan Li, and Dongsheng Li. 2024b. Llms can find mathematical reasoning mistakes by pedagogical chain-of-thought . In <i>Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI 2024, Jeju, South Korea, August 3-9, 2024</i> , pages 3439–3447. ijcai.org.	1130
1079			1131
1080			1132
1081			1133
1082			1134
1083			
1084	Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Live-codebench: Holistic and contamination free evaluation of large language models for code . <i>ArXiv preprint</i> , abs/2403.07974.	Mingyu Jin, Qinkai Yu, Dong Shu, Haiyan Zhao, Wenyue Hua, Yanda Meng, Yongfeng Zhang, and Mengnan Du. 2024. The impact of reasoning step length on large language models . <i>ArXiv preprint</i> , abs/2401.04925.	1135
1085			1136
1086			1137
1087			1138
1088			1139
1089			
1090	Joel Jang, Seonghyeon Ye, and Minjoon Seo. 2023. Can large language models truly understand prompts? a	Muhammad Khalifa, Lajanugen Logeswaran, Moontae Lee, Honglak Lee, and Lu Wang. 2023. GRACE: Discriminator-guided chain-of-thought reasoning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 15299–15328, Singapore. Association for Computational Linguistics.	1140
1091			1141
			1142
			1143
			1144
			1145

1146	Heegy Kim, Taeyang Jeon, Seunghwan Choi, Seungtaek Choi, and Hyunsouk Cho. 2024a. Flex: Expert-level false-less execution metric for reliable text-to-sql benchmark . <i>ArXiv preprint</i> , abs/2409.19014.	language models and chain-of-thought for automatic scoring. <i>Computers and Education: Artificial Intelligence</i> , 6:100213.	1203 1204 1205
1150	Seungone Kim, Jamin Shin, Yejin Choi, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2024b. Prometheus: Inducing fine-grained evaluation capability in language models . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	Jaehyeok Lee, Keisuke Sakaguchi, and JinYeong Bak. 2025a. Self-training meets consistency: Improving llms' reasoning with consistency-driven rationale evaluation. In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 10519–10539.	1206 1207 1208 1209 1210 1211 1212 1213
1158	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	Jinu Lee and Julia Hockenmaier. 2025. Evaluating step-by-step reasoning traces: A survey . <i>ArXiv preprint</i> , abs/2502.12289.	1214 1215 1216
1165	Aakriti Kumar, Nalin Pounpeth, Diyi Yang, Erina Farrell, Bruce Lambert, and Matthew Groh. 2025. When large language models are reliable for judging empathic communication . <i>ArXiv preprint</i> , abs/2506.10150.	Young-Jun Lee, Seungone Kim, Byung-Kwan Lee, Minkyong Moon, Yechan Hwang, Jong Myoung Kim, Graham Neubig, Sean Welleck, and Ho-Jin Choi. 2025b. Refinebench: Evaluating refinement capability of language models via checklists . <i>arXiv preprint arXiv:2511.22173</i> .	1217 1218 1219 1220 1221 1222
1170	Lucas Fonseca Lage and Simon Ostermann. 2025. Openfactscore: Open-source atomic evaluation of factuality in text generation . <i>ArXiv preprint</i> , abs/2507.05965.	Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay V. Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. Solving quantitative reasoning problems with language models . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	1223 1224 1225 1226 1227 1228 1229 1230 1231 1232
1174	Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-Tau Yih, Daniel Fried, Sida I. Wang, and Tao Yu. 2023. DS-1000: A natural and reliable benchmark for data science code generation . 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 18319–18345. PMLR.	Bohan Li, Jiannan Guan, Longxu Dou, Yunlong Feng, Dingzirui Wang, Yang Xu, Enbo Wang, Qiguang Chen, Bichen Wang, Xiao Xu, and 1 others. 2025a. Can large language models understand you better? an mbti personality detection dataset aligned with population traits . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 5071–5081.	1233 1234 1235 1236 1237 1238 1239 1240
1183	Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, and 1 others. 2023. Measuring faithfulness in chain-of-thought reasoning . <i>ArXiv preprint</i> , abs/2307.13702.	Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, and 1 others. 2025b. From generation to judgment: Opportunities and challenges of llm-as-a-judge . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 2757–2791.	1241 1242 1243 1244 1245 1246 1247 1248
1189	Jon M Laurent, Joseph D Janizek, Michael Ruzo, Michaela M Hinks, Michael J Hammerling, Siddharth Narayanan, Manvitha Ponnampati, Andrew D White, and Samuel G Rodrigues. 2024. Lab-bench: Measuring capabilities of language models for biology research . <i>ArXiv preprint</i> , abs/2407.10362.	Jia Li, Ge Li, Yongmin Li, and Zhi Jin. 2025c. Structured chain-of-thought prompting for code generation . <i>ACM Transactions on Software Engineering and Methodology</i> , 34(2):1–23.	1249 1250 1251 1252
1195	Joshua Ong Jun Leang, Aryo Pradipta Gema, and Shay B Cohen. 2025. Comat: Chain of mathematically annotated thought improves mathematical reasoning . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 20256–20285.	Junyi Li, Xiaoxue Cheng, Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023a. HaluEval: A large-scale hallucination evaluation benchmark for large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 6449–6464, Singapore. Association for Computational Linguistics.	1253 1254 1255 1256 1257 1258 1259
1201	Gyeong-Geon Lee, Ehsan Latif, Xuansheng Wu, Ninghao Liu, and Xiaoming Zhai. 2024. Applying large		

1260	Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, and Deheng Ye. 2024a. More agents is all you need. <i>arXiv preprint arXiv:2402.05120</i> .		
1261		Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang, Mingu Lee, Roland Memisevic, and Hao Su. 2023a. Deductive verification of chain-of-thought reasoning . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1313
1262			1314
1263	Ruosen Li and Xinya Du. 2023. Leveraging structured information for explainable multi-hop question answering and reasoning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 6779–6789, Singapore. Association for Computational Linguistics.		1315
1264			1316
1265			1317
1266			1318
1267		Zhan Ling and 1 others. 2023b. Thought propagation: An analogical approach to complex reasoning . volume abs/2310.03965.	1319
1268			1320
1269	Ruosen Li, Ziming Luo, and Xinya Du. 2024b. Fine-grained hallucination detection and mitigation in language model mathematical reasoning.	Chengwu Liu, Ye Yuan, Yichun Yin, Yan Xu, Xin Xu, Zaoyu Chen, Yasheng Wang, Lifeng Shang, Qun Liu, and Ming Zhang. 2025a. Safe: Enhancing mathematical reasoning in large language models via retrospective step-aware formal verification . <i>ArXiv preprint, abs/2506.04592</i> .	1323
1270			1324
1271			1325
1272	Xiaoyuan Li, Wenjie Wang, Moxin Li, Junrong Guo, Yang Zhang, and Fuli Feng. 2024c. Evaluating mathematical reasoning of large language models: A focus on error identification and correction . <i>ArXiv preprint, abs/2406.00755</i> .		1326
1273			1327
1274			1328
1275		Hanmeng Liu, Zhiyang Teng, Leyang Cui, Chaoli Zhang, Qiji Zhou, and Yue Zhang. 2023a. Logi-CoT: Logical chain-of-thought instruction tuning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 2908–2921, Singapore. Association for Computational Linguistics.	1329
1276			1330
1277	Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2023b. Diverse: A reliable approach to few-shot reasoning. In <i>ACL</i> .		1331
1278			1332
1279			1333
1280	Yiwei Li, Peiwen Yuan, Shaoxiong Feng, Boyuan Pan, Xinglin Wang, Bin Sun, Heda Wang, and Kan Li. 2024d. Escape sky-high cost: Early-stopping self-consistency for multi-step reasoning . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	Hugo Liu and Push Singh. 2004. Conceptnet—a practical commonsense reasoning tool-kit. <i>BT technology journal</i> , 22(4):211–226.	1334
1281			1335
1282			1336
1283			1337
1284		Jiacheng Liu, Alisa Liu, Ximing Lu, Sean Welleck, Peter West, Ronan Le Bras, Yejin Choi, and Hannaneh Hajishirzi. 2022. Generated knowledge prompting for commonsense reasoning . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3154–3169, Dublin, Ireland. Association for Computational Linguistics.	1338
1285			1339
1286			1340
1287	Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, and 1 others. 2022. Competition-level code generation with alphacode. <i>Science</i> , 378(6624):1092–1097.		1341
1288			1342
1289			1343
1290			1344
1291			1345
1292	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, and 1 others. 2022. Holistic evaluation of language models . <i>ArXiv preprint, abs/2211.09110</i> .	Tianqiao Liu, Zui Chen, Zitao Liu, Mi Tian, and Weiqi Luo. 2024. Expediting and elevating large language model reasoning via hidden chain-of-thought decoding . <i>ArXiv preprint, abs/2409.08561</i> .	1346
1293			1347
1294			1348
1295			1349
1296		Tongxuan Liu, Wenjiang Xu, Weizhe Huang, Yuting Zeng, Jiaxing Wang, Xingyu Wang, Hailong Yang, and Jing Li. 2025b. Logic-of-thought: Injecting logic into contexts for full reasoning in large language models. In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 10168–10185.	1350
1297	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Rui Wang, Zhaopeng Tu, Shuming Shi, and Xing Wang. 2024. Encouraging divergent thinking in large language models through multi-agent debate. In <i>EMNLP</i> .		1351
1298			1352
1299			1353
1300			1354
1301			1355
1302	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. Let’s verify step by step . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023b. G-eval: NLG evaluation using gpt-4 with better human alignment . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics.	1356
1303			1357
1304			1358
1305			1359
1306			1360
1307			1361
1308			1362
1309	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, and 1 others. 2021. Codexglue: A machine learning benchmark	1363
1310			1364
1311			1365
1312			1366

1370		dataset for code understanding and generation. <i>ArXiv preprint</i> , abs/2102.04664.	
1371			
1372	Scott M. Lundberg and Su-In Lee. 2017.	A unified approach to interpreting model predictions. In <i>Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA</i> , pages 4765–4774.	
1373			
1374			
1375			
1376			
1377			
1378	Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. 2024.	Reasoning on graphs: Faithful and interpretable large language model reasoning. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	
1379			
1380			
1381			
1382			
1383			
1384	Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, D Deligiannis, and 1 others. 2024.	Walk the talk? measuring the faithfulness of large language model explanations. In <i>ICLR</i> .	
1385			
1386			
1387			
1388	Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023a.	Faithful chain-of-thought reasoning. In <i>Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 305–329, Nusa Dua, Bali. Association for Computational Linguistics.	
1389			
1390			
1391			
1392			
1393			
1394			
1395			
1396			
1397			
1398	Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023b.	Faithful chain-of-thought reasoning. In <i>Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 305–329, Nusa Dua, Bali. Association for Computational Linguistics.	
1399			
1400			
1401			
1402			
1403			
1404			
1405			
1406			
1407			
1408	Aman Madaan, Katherine Hermann, and Amir Yazdanbakhsh. 2023a.	What makes chain-of-thought prompting effective? a counterfactual study. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 1448–1535, Singapore. Association for Computational Linguistics.	
1409			
1410			
1411			
1412			
1413			
1414	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023b.	Self-refine: Iterative refinement with self-feedback. In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	
1415			
1416			
1417			
1418			
1419			
1420			
1421			
1422			
1423			
1424			
1425	Bodhisattwa Prasad Majumder, Harshit Surana, Dhruv Agarwal, Bhavana Dalvi Mishra, Abhijeetsingh		
1426			
	Meena, Aryan Prakhar, Tirth Vora, Tushar Khot, Ashish Sabharwal, and Peter Clark. 2024.	Discoverybench: Towards data-driven discovery with large language models. <i>ArXiv preprint</i> , abs/2407.01725.	1427 1428 1429 1430
	Youssef Maklad, Fares Wael, Wael Elersy, and Ali Hamdi. 2025.	Retrieval augmented generation based llm evaluation for protocol state machine inference with chain-of-thought reasoning. In <i>International Congress on Information and Communication Technology</i> , pages 313–323. Springer.	1431 1432 1433 1434 1435 1436
	Zhenjiang Mao, Artem Bisliouk, Rohith Reddy Nama, and Ivan Ruchkin. 2025.	Temporalizing confidence: Evaluation of chain-of-thought reasoning with signal temporal logic. <i>ArXiv preprint</i> , abs/2506.08243.	1437 1438 1439 1440
	Ning Miao, Yee Whye Teh, and Tom Rainforth. 2024.	Selfcheck: Using llms to zero-shot check their own step-by-step reasoning. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	1441 1442 1443 1444 1445
	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018.	Can a suit of armor conduct electricity? a new dataset for open book question answering. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2381–2391, Brussels, Belgium. Association for Computational Linguistics.	1446 1447 1448 1449 1450 1451 1452
	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023.	FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12076–12100, Singapore. Association for Computational Linguistics.	1453 1454 1455 1456 1457 1458 1459 1460
	Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022.	Rethinking the role of demonstrations: What makes in-context learning work? In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11048–11064, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1461 1462 1463 1464 1465 1466 1467 1468
	Stephen Miner, Yoshiki Takashima, Simeng Han, Sam Kouteili, Ferhat Erata, Ruzica Piskac, and Scott J Shapiro. 2024.	Scheherazade: Evaluating chain-of-thought math reasoning in llms with chain-of-problems. <i>ArXiv preprint</i> , abs/2410.00151.	1469 1470 1471 1472 1473
	Francesco Maria Molfese, Luca Moroni, Ciro Porcaro, Simone Conia, and Roberto Navigli. 2025.	Re-traceqa: Evaluating reasoning traces of small language models in commonsense question answering. <i>ArXiv preprint</i> , abs/2510.09351.	1474 1475 1476 1477 1478
	Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori B Hashimoto. 2025.	s1: Simple test-time scaling. In <i>Proceedings of the 2025 Conference on</i>	1479 1480 1481 1482 1483

1484	<i>Empirical Methods in Natural Language Processing</i> ,	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	1538
1485	pages 20286–20332.	Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	1539
1486	S Mysore and 1 others. 2023. Crisp: Complex reasoning		1540
1487	with interpretable step-based plans. In <i>EMNLP</i> .		1541
1488	Sina Bagheri Nezhad, Yao Li, and Ameeta Agrawal.		1542
1489	2025. Symcode: A neurosymbolic approach to mathematical reasoning via verifiable code generation .		1543
1490	<i>ArXiv preprint</i> , abs/2510.25975.	Arkil Patel, Satwik Bhattamishra, and Navin Goyal.	1544
1491		2021. Are NLP models really able to solve simple math word problems? In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2080–2094, Online. Association for Computational Linguistics.	1545
1492	Minh-Vuong Nguyen, Linhao Luo, Fatemeh Shiri, Dinh		1546
1493	Phung, Yuan-Fang Li, Thuy Vu, and Gholamreza		1547
1494	Haffari. 2024. Direct evaluation of chain-of-thought		1548
1495	in multi-hop reasoning with knowledge graphs. In		1549
1496	<i>Findings of the Association for Computational Lin-</i>		1550
1497	<i>guistics: ACL 2024</i> , pages 2862–2883.		1551
1498	Phuong Minh Nguyen, Tien Huu Dang, and Naoya In-		
1499	oue. 2025. Non-interactive symbolic-aided chain-		
1500	of-thought for logical reasoning . <i>ArXiv preprint</i> ,		
1501	abs/2508.12425.	Kaviraj Pather, Elena Hadjigeorgiou, Arben Krasniqi,	1552
1502		Claire Schmit, Irina Rusu, Marc Pons, and Kabir	1553
1503		Khan. 2025. Vis-cot: A human-in-the-loop frame-	1554
1504		work for interactive visualization and intervention	1555
1505		in llm chain-of-thought reasoning . <i>ArXiv preprint</i> ,	1556
1506		abs/2509.01412.	1557
1507			
1508	Natapong Nitarach, Warit Sirichotedumrong, Panop		
1509	Pitchayarthorn, Pittawat Taveekitworachai, Potsawee		
1510	Manakul, and Kunat Pipatanakul. 2025. Fincot:		
1511	Grounding chain-of-thought in expert financial rea-		
1512	soning . <i>ArXiv preprint</i> , abs/2506.16123.		
1513		Debjit Paul, Robert West, Antoine Bosselut, and Boi	1558
1514		Faltings. 2024. Making reasoning matter: Measur-	1559
1515		ing and improving faithfulness of chain-of-thought	1560
1516		reasoning . <i>ArXiv preprint</i> , abs/2402.13950.	1561
1517			
1518		Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng,	1562
1519		Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou	1563
1520		Yu, Weizhu Chen, and 1 others. 2023. Check your	1564
1521		facts and try again: Improving large language models	1565
1522		with external knowledge and automated feedback .	1566
1523		<i>ArXiv preprint</i> , abs/2302.12813.	1567
1524			
1525		Gabriel Poesia, Kanishk Gandhi, Eric Zelikman, and	1568
1526		Noah D Goodman. 2023. Certified deductive rea-	1569
1527		soning with language models . <i>ArXiv preprint</i> ,	1570
1528		abs/2306.04031.	1571
1529			
1530		Nearchos Potamitis, Lars Klein, and Akhil Arora. 2025.	1572
1531		Reasonbench: Benchmarking the (in) stability of llm	1573
1532		reasoning . <i>arXiv preprint arXiv:2512.07795</i> .	1574
1533			
1534		Archiki Prasad, Swarnadeep Saha, Xiang Zhou, and	1575
1535		Mohit Bansal. 2023. ReCEval: Evaluating reasoning	1576
1536		chains via correctness and informativeness . In <i>Pro-</i>	1577
1537		<i>ceedings of the 2023 Conference on Empirical Meth-</i>	1578
		<i>ods in Natural Language Processing</i> , pages 10066–	1579
		10086, Singapore. Association for Computational	1580
		Linguistics.	1581
		Yukun Qi, Yiming Zhao, Yu Zeng, Xikun Bao, Wenxuan	1582
		Huang, Lin Chen, Zehui Chen, Jie Zhao, Zhongang	1583
		Qi, and Feng Zhao. 2025. Vcr-bench: A compre-	1584
		hensive evaluation framework for video chain-of-thought	1585
		reasoning . <i>ArXiv preprint</i> , abs/2504.07956.	1586
		Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Kunlun	1587
		Zhu, Hanchen Xia, Yufan Dang, Zhuoyun Du, Weize	1588
		Chen, Cheng Yang, and 1 others. 2024. Scaling	1589
		large language model-based multi-agent collabora-	1590
		tion . <i>arXiv preprint arXiv:2406.07155</i> .	1591
		Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen,	1592
		Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang,	1593

1822	Miles Turpin, Julian Michael, Ethan Perez, and Samuel R. Bowman. 2023. Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	<i>the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4958–4981.	1878
1823			1879
1824			1880
1825			
1826		Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, and 1 others. 2024b. A survey on large language model based autonomous agents. <i>Frontiers of Computer Science</i> , 18(6):186345.	1881
1827			1882
1828			1883
1829			1884
			1885
1830	Martin Tutek, Fateme Hashemi Chaleshtori, Ana Marasović, and Yonatan Belinkov. 2025. Measuring faithfulness of chains of thought by unlearning reasoning steps. <i>arXiv e-prints</i> , pages arXiv–2502.	Linhao Wang, Zihan Wang, Jinghao Lin, Qiwei Zeng, Jingdan Zhang, Tianrun Wu, Xuanhao Zhang, Junjie Li, Zijing Wang, Yongkang Liu, and 1 others. 2025c. Good teachers, better students: A survey of reward models for llm. <i>Authorea Preprints</i> .	1886
1831			1887
1832			1888
1833			1889
1834	Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. 2022. Solving math word problems with process-and outcome-based feedback . <i>ArXiv preprint</i> , abs/2211.14275.		1890
1835			
1836		Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. 2024c. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 9426–9439.	1891
1837			1892
1838			1893
1839	Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. 2022. Large language models still can't plan (a benchmark for llms on planning and reasoning about change). In <i>NeurIPS 2022 Foundation Models for Decision Making Workshop</i> .		1894
1840			1895
1841			1896
1842			1897
1843		Tianyu Wang, Nianjun Zhou, and Zhixiong Chen. 2024d. Enhancing computer programming education with llms: A study on effective prompt engineering for python code generation . <i>ArXiv preprint</i> , abs/2407.05437.	1898
1844	Hemish Veeraboina. 2023. Aime problem set 1983-2024 .		1899
1845			1900
1846			1901
1847	Theo Walker, Christopher M Grulke, Diane Pozefsky, and Alexander Tropsha. 2010. Chembench: a cheminformatics workbench. <i>Bioinformatics</i> , 26(23):3000–3001.		1902
1848			
1849			
1850	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding . In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 353–355, Brussels, Belgium. Association for Computational Linguistics.	Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R. Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2024e. Scibench: Evaluating college-level scientific problem-solving abilities of large language models . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	1903
1851			1904
1852			1905
1853			1906
1854			1907
1855			1908
1856			1909
1857			1910
1858	Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. 2023a. Towards understanding chain-of-thought prompting: An empirical study of what matters . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2717–2739, Toronto, Canada. Association for Computational Linguistics.	Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2024f. MINT: evaluating llms in multi-turn interaction with tools and language feedback . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	1911
1859			1912
1860			1913
1861			1914
1862			1915
1863			1916
1864		Xinyuan Wang and 1 others. 2023b. Is reasoning a hindsight? evaluating the faithfulness of chain-of-thought explanations. In <i>EMNLP</i> .	1917
1865			1918
1866	Boxuan Wang, Zhuoyun Li, Xinmiao Huang, Xiaowei Huang, and Yi Dong. 2025a. Chasing consistency: Quantifying and optimizing human-model alignment in chain-of-thought reasoning . <i>ArXiv preprint</i> , abs/2511.06168.		1919
1867			
1868		Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023c. 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.	1920
1869			1921
1870			1922
1871	Changyue Wang, Weihang Su, Qingyao Ai, and Yiqun Liu. 2025b. Joint evaluation of answer and reasoning consistency for hallucination detection in large reasoning models . <i>ArXiv preprint</i> , abs/2506.04832.		1923
1872			1924
1873			1925
1874			1926
1875	Jianing Wang, Qiushi Sun, Xiang Li, and Ming Gao. 2024a. Boosting language models reasoning with chain-of-knowledge prompting. In <i>Proceedings of</i>	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023d. 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.	1927
1876			1928
1877			1929
			1930
			1931
			1932
			1933

2043	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafraan, Karthik R. Narasimhan, and Yuan Cao. 2023b. React: Synergizing reasoning and acting in language models . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	2098
2044		2099
2045		2100
2046		2101
2047		2102
2048		
2049	Xi Ye, Qiaochu Chen, Isil Dillig, and Greg Durrett. 2023a. Satlm: Satisfiability-aided language models using declarative prompting . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	2103
2050		2104
2051		2105
2052		2106
2053		2107
2054		
2055		
2056	Xi Ye and 1 others. 2023b. A comprehensive study of chain-of-thought reasoning in large language models: Advances, challenges, and opportunities . <i>ArXiv preprint</i> , abs/2309.15402.	2108
2057		2109
2058		2110
2059		2111
2060	Wenhan Yu, Xinbo Lin, Lanxin Ni, Jinhua Cheng, and Lei Sha. 2025. Benchmarking multi-step legal reasoning and analyzing chain-of-thought effects in large language models . <i>ArXiv preprint</i> , abs/2511.07979.	2112
2061		2113
2062		
2063		
2064		
2065	Zihan Yu, Liang He, Zhen Wu, Xinyu Dai, and Jiajun Chen. 2023. Towards better chain-of-thought prompting strategies: A survey . <i>ArXiv preprint</i> , abs/2310.04959.	2114
2066		2115
2067		2116
2068		2117
2069	Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. Bartscore: Evaluating generated text as text generation . In <i>Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual</i> , pages 27263–27277.	2118
2070		2119
2071		
2072		
2073		
2074		
2075	Weiqi Yue, Yuyu Yin, Xin Zhang, Binbin Shi, Tingting Liang, and Jian Wan. 2025. Cot4rec: Revealing user preferences through chain of thought for recommender systems . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 13142–13151.	2120
2076		2121
2077		2122
2078		2123
2079		2124
2080		2125
2081	Zhongshen Zeng, Pengguang Chen, Shu Liu, Haiyun Jiang, and Jiaya Jia. 2023. Mr-gsm8k: A meta-reasoning benchmark for large language model evaluation . <i>ArXiv preprint</i> , abs/2312.17080.	2126
2082		
2083		
2084		
2085	Bingyuan Zhang, Xulong Zhang, Yong Zhang, Jun Yu, and Jianzong Wang. 2025a. Logic consistency makes large language models personalized reasoning teachers . In <i>2025 International Joint Conference on Neural Networks (IJCNN)</i> , pages 1–8. IEEE.	2127
2086		2128
2087		2129
2088		2130
2089		2131
2090	Boning Zhang, Chengxi Li, and Kai Fan. 2024a. Mario eval: Evaluate your math llm with your math llm—a mathematical dataset evaluation toolkit . <i>ArXiv preprint</i> , abs/2404.13925.	2132
2091		2133
2092		2134
2093		2135
2094	Jieyu Zhang, Ranjay Krishna, Ahmed H Awadallah, and Chi Wang. 2023a. Ecoassistant: Using llm assistant more affordably and accurately . <i>ArXiv preprint</i> , abs/2310.03046.	2136
2095		2137
2096		
2097		
	Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. 2025b. Generative verifiers: Reward modeling as next-token prediction . In <i>The Thirteenth International Conference on Learning Representations</i> .	2140
		2141
		2142
		2143
	Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A. Smith. 2024b. How language model hallucinations can snowball . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	2144
		2145
		2146
		2147
		2148
	Qiyuan Zhang, Fuyuan Lyu, Zexu Sun, Lei Wang, Weixu Zhang, Wenyue Hua, Haolun Wu, Zhihan Guo, Yufei Wang, Niklas Muennighoff, and 1 others. 2025c. A survey on test-time scaling in large language models: What, how, where, and how well? <i>arXiv preprint arXiv:2503.24235</i> .	2149
		2150
		2151
		2152
	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with BERT . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	
	Xuan Zhang, Chao Du, Tianyu Pang, Qian Liu, Wei Gao, and Min Lin. 2024c. Chain of preference optimization: Improving chain-of-thought reasoning in llms . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	
	Zhenru Zhang, Chujie Zheng, Yangzhen Wu, Beichen Zhang, Runji Lin, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. 2025d. The lessons of developing process reward models in mathematical reasoning . <i>ArXiv preprint</i> , abs/2501.07301.	
	Zhuosheng Zhang, Yao Yao, Aston Zhang, Xiangru Tang, Xinbei Ma, Zhiwei He, Yiming Wang, Mark Gerstein, Rui Wang, Gongshen Liu, and 1 others. 2025e. Igniting language intelligence: The hitchhiker’s guide from chain-of-thought reasoning to language agents . <i>ACM Computing Surveys</i> , 57(8):1–39.	
	Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2023b. Automatic chain of thought prompting in large language models . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
	Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2024. Explainability for large language models: A survey . <i>ACM Transactions on Intelligent Systems and Technology</i> , 15(2):1–38.	
	Kesen Zhao, Beier Zhu, Qianru Sun, and Hanwang Zhang. 2025a. Unsupervised visual chain-of-thought reasoning via preference optimization . <i>arXiv preprint arXiv:2504.18397</i> .	

- 2153 Ruochen Zhao, Xingxuan Li, Shafiq Joty, Chengwei
2154 Qin, and Lidong Bing. 2023. [Verify-and-edit: A](#)
2155 [knowledge-enhanced chain-of-thought framework](#).
2156 In *Proceedings of the 61st Annual Meeting of the*
2157 *Association for Computational Linguistics (Volume*
2158 *1: Long Papers)*, pages 5823–5840, Toronto, Canada.
2159 Association for Computational Linguistics.
- 2160 Zheng Zhao, Yeskendir Koishakenov, Xianjun Yang,
2161 Naila Murray, and Nicola Cancedda. 2025b. [Verify-](#)
2162 [ing chain-of-thought reasoning via its computational](#)
2163 [graph](#). *ArXiv preprint*, abs/2510.09312.
- 2164 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
2165 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
2166 Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,
2167 Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging](#)
2168 [llm-as-a-judge with mt-bench and chatbot arena](#). In
2169 *Advances in Neural Information Processing Systems*
2170 *36: Annual Conference on Neural Information Pro-*
2171 *cessing Systems 2023, NeurIPS 2023, New Orleans,*
2172 *LA, USA, December 10 - 16, 2023*.
- 2173 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei,
2174 Nathan Scales, Xuezhi Wang, Dale Schuurmans,
2175 Claire Cui, Olivier Bousquet, Quoc V. Le, and Ed H.
2176 Chi. 2023. [Least-to-most prompting enables com-](#)
2177 [plex reasoning in large language models](#). In *The*
2178 *Eleventh International Conference on Learning Rep-*
2179 *resentations, ICLR 2023, Kigali, Rwanda, May 1-5,*
2180 *2023*. OpenReview.net.
- 2181 Runtao Zhou, Giang Nguyen, Nikita Kharya, Anh
2182 Nguyen, and Chirag Agarwal. 2025. [Improving](#)
2183 [human verification of llm reasoning through in-](#)
2184 [teractive explanation interfaces](#). *ArXiv preprint*,
2185 abs/2510.22922.
- 2186 Zhanke Zhou, Rong Tao, Jianing Zhu, Yiwen Luo, Zeng-
2187 ma Wang, and Bo Han. 2024. [Can language models](#)
2188 [perform robust reasoning in chain-of-thought prompt-](#)
2189 [ing with noisy rationales?](#) In *Advances in Neural*
2190 *Information Processing Systems 38: Annual Confer-*
2191 *ence on Neural Information Processing Systems 2024,*
2192 *NeurIPS 2024, Vancouver, BC, Canada, December*
2193 *10 - 15, 2024*.
- 2194 Kaijie Zhu, Qinlin Zhao, Hao Chen, Jindong Wang, and
2195 Xing Xie. 2024. Promptbench: A unified library
2196 for evaluation of large language models. *Journal of*
2197 *Machine Learning Research*, 25(254):1–22.
- 2198 Mingchen Zhuge, Haozhe Liu, Francesco Faccio, Dy-
2199 lan R Ashley, Róbert Csordás, Anand Gopalakrish-
2200 nan, Abdullah Hamdi, Hasan Abed Al Kader Ham-
2201 moud, Vincent Herrmann, Kazuki Irie, and 1 others.
2202 2023. Mindstorms in natural language-based soci-
2203 eties of mind. *arXiv preprint arXiv:2305.17066*.

A Example of Disguised Accuracy

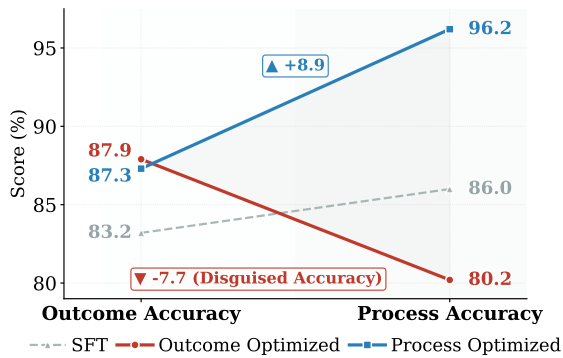


Figure 5: An example showing how high answer accuracy can mask poor intermediate reasoning.

Figure 5 illustrates experimental results from (Uesato et al., 2022) using a 70B model trained with Supervised Fine-Tuning (SFT), Outcome-Optimized, and Process-Optimized methods, showing the disguised accuracy phenomenon (Bentham et al., 2024). In these experiments, both the answer accuracy and the accuracy of the reasoning process are evaluated on the same dataset. The results show that while the Outcome-Optimized method improves final outcome accuracy, it causes a severe degradation in the quality of intermediate reasoning steps. This indicates that the model guesses correct answers using flawed logic rather than sound reasoning. As a result, this phenomenon leads to a reversal in model rankings, where a model with higher outcome metrics actually demonstrates inferior reasoning capabilities. In practical applications, such models provide low trustworthiness because their high accuracy masks fundamental reasoning failures.

B Related Work

Surveys on CoT Generation and Prompting. A substantial body of literature reviews the landscape of CoT reasoning, primarily focusing on generation strategies and prompt engineering (Yu et al., 2023; Chu et al., 2024; Zhang et al., 2025e; Qiao et al., 2023). These surveys extensively catalog paradigms ranging from Zero-shot (Kojima et al., 2022) and Few-shot CoT (Wei et al., 2022) to complex agentic frameworks (Chen et al., 2025c; Xia et al., 2025). However, their analytical lens is predominantly cast on *how to elicit* reasoning rather than *how to evaluate* it. Consequently, assessment in these works is often reduced to outcome-oriented

metrics (e.g., accuracy on GSM8K), largely conflating the validity of the intermediate process with the correctness of the final answer (Lee and Hockenmaier, 2025).

General LLM Evaluation and Specific Sub-domains. Broader reviews of LLM evaluation (Chang et al., 2024; Guo et al., 2023; Liang et al., 2022; Srivastava et al., 2023) treat reasoning as merely one of many capabilities. While they establish benchmarks for general performance, they rarely dissect the granular quality of reasoning traces, such as redundancy or step-wise necessity. Similarly, while surveys on LLM-as-a-Judge (Wang et al., 2024b; Zheng et al., 2023) and hallucination (Ji et al., 2023a; Huang et al., 2025a) touch upon evaluation mechanics, they typically focus on linguistic fluency, preference alignment, or factual fabrication, lacking a systematic taxonomy for the causal rigor and logical consistency inherent to multi-step reasoning.

To date, there is no comprehensive survey dedicated specifically to CoT evaluation that systematically bridges theoretical dimensions with practical methodologies across diverse domains. Our work addresses this critical void by establishing a unified framework structured along four core dimensions: **What** to measure, **How** to verify, **Who** evaluates, and **Where** to apply.

C Integration of Interpretability Methods

To assess CoT fidelity, recent studies integrate the explainability methods. Researchers utilize tools like LIME (Ribeiro et al., 2016) and SHAP (Lundberg and Lee, 2017), alongside gradient-based self-influence methods (Thakkar et al., 2023), to cross-verify if implicit token importance aligns with explicit reasoning steps. Counterfactual perturbations test causal validity by determining whether the generated rationale consistently dictates the prediction shift under input changes (Turpin et al., 2023; Atanasova et al., 2023; Wang et al., 2023b). This combined framework helps distinguish authentic reasoning processes from persuasive yet unfaithful post-hoc rationalizations (Lyu et al., 2024; Lanham et al., 2023).

D Benchmarks & Evaluation Methodologies

Table 1 provides an overview of the benchmark datasets across different tasks mentioned in this survey for CoT evaluation.

2288	D.1 Mathematical Reasoning		
2289	D.1.1 Benchmarks		
2290	Mathematical evaluation prioritizes absolute Correctness , yet complex problems demand higher		2337
2291	Logical Coherence . Basic robustness relies on		2338
2292	GSM8K (Cobbe et al., 2021) and SVAMP (Patel		2339
2293	et al., 2021). Conversely, MathQA (Amini et al.,		2340
2294	2019) and MATH (Hendrycks et al., 2021) target		2341
2295	interpretable skills where Cognitive Efficiency of-		2342
2296	ten yields to detailed derivation. MATH further		2343
2297	includes challenging problems like AMC12 and		2344
2298	AIME (Veeraboina, 2023). To address saturation,		
2299	recent works increase complexity to test Logical		
2300	Coherence . GSM8K-Scheherazade (Miner et al.,		
2301	2024) uses chained queries, while CoMAT (Leang		
2302	et al., 2025) and SCoT (Wang et al., 2024h) enforce		
2303	strategy elicitation.		
2304			
2305	D.1.2 Evaluation Methodologies		
2306	Automated Verification. Automated strategies		
2307	split into symbolic verification for Correct-		
2308	ness and semantic scoring for Interpretability .		
2309	Tools like MathSensei (Das et al., 2024), Sym-		
2310	Code (Nezhad et al., 2025), and RLSR (Simonds		
2311	et al., 2025) leverage SymPy for execution, while		
2312	Safe (Liu et al., 2025a) uses retrospective verifi-		
2313	cation. To assess Logical Coherence , verifiers		
2314	from GSM8K (Cobbe et al., 2021) evolve into		
2315	step-wise discriminators in CoMAT (Leang et al.,		
2316	2025) and Hidden CoT (Liu et al., 2024). Addition-		
2317	ally, Nguyen et al. (2024) use knowledge graphs		
2318	for semantic alignment, aggregated by MARIO		
2319	Eval (Zhang et al., 2024a).		
2320	Human-in-the-Loop. Human evaluation re-		
2321	remains the ground truth. Experts are essential		
2322	for assessing Logical Coherence in complex		
2323	proofs (Zeng et al., 2023) and providing process su-		
2324	pervision (Lightman et al., 2024), serving as a gold		
2325	standard (Nguyen et al., 2024). Conversely, crowd-		
2326	sourcing offers scalability for subjective Inter-		
2327	pretability in ROSCOE (Golovneva et al., 2023),		
2328	despite limitations in complex logic. Interactive		
2329	interfaces (Zhou et al., 2025) help mitigate this by		
2330	visualizing reasoning chains for non-experts.		
2331	D.2 Commonsense Reasoning		
2332	D.2.1 Benchmarks		
2333	Commonsense domains prioritize Interpretability		
2334	and Cognitive Efficiency . Benchmarks include		
2335	CommonsenseQA (Talmor et al., 2019), Strate-		
2336	gyQA (Geva et al., 2021), ARC (Clark et al.,		
	2018), and OpenBookQA (Mihaylov et al., 2018).		2337
	ROSCOE (Golovneva et al., 2023) quantifies se-		2338
	matic Consistency , while SELF-CHECK (Miao		2339
	et al., 2024) detects errors. To ensure Correctness ,		2340
	SciKnowEval (Feng et al., 2024) and CoK (Wang		2341
	et al., 2024a) integrate knowledge triples, while De-		2342
	CoT (Wu et al., 2024b) mitigates bias. VISCO (Wu		2343
	et al., 2025a) extends this to multimodal tasks.		2344
	D.2.2 Evaluation Methodologies		2345
	Hybrid Evaluation. Hybrid methods are com-		2346
	mon, using frameworks like RORA (Jiang et al.,		2347
	2024a) and ACORN (Brassard et al., 2024) to align		2348
	rationales with human Interpretability . Along		2349
	with ReTraceQA (Molfese et al., 2025), these		2350
	show that open-weight models often need specific		2351
	prompting to match human intuition (Zheng et al.,		2352
	2023; Bang et al., 2023).		2353
	Analysis of Pitfall. Methodologies must diag-		2354
	nose violations of Cognitive Efficiency , such as		2355
	bias or over-reasoning. CoT can amplify stereo-		2356
	types (Shaikh et al., 2023), while over-explaining		2357
	axioms leads to hallucinations. Protocols test for		2358
	unfaithful explanations that fail to support predic-		2359
	tions (Turpin et al., 2023).		2360
	D.3 Code Reasoning		2361
	D.3.1 Benchmarks		2362
	Code reasoning demands strict Correctness via		2363
	execution while using Logical Coherence for plan-		2364
	ning. Benchmarks range from HumanEval (Chen,		2365
	2021) and MBPP (Austin et al., 2021) to com-		2366
	petition sets like APPS (Hendrycks et al., 2021)		2367
	and CodeContests (Li et al., 2022). Multilingual		2368
	tests include MBXP and MathQA-X (Athiwaratkun		2369
	et al., 2023). Specialized domains cover data sci-		2370
	ence in DS-1000 (Lai et al., 2023) and general		2371
	understanding in CodeXGLUE (Lu et al., 2021).		2372
	Recent tools like CRUXEval (Gu et al., 2024) and		2373
	LiveCodeBench (Jain et al., 2024) expand this to		2374
	execution reasoning.		2375
	D.3.2 Evaluation Methodologies		2376
	Automated Evaluation. Protocols assess func-		2377
	tional correctness via execution metrics like		2378
	pass@k (Chen, 2021). Complementing this, static		2379
	analysis tools like Pylint detect syntactic viola-		2380
	tions and enforce style Consistency , identifying		2381
	latent errors without execution (Salih and Sarhan,		2382
	2025; Jiang et al., 2025c; Almeida et al., 2024;		2383
	Wang et al., 2024d; Al-Khafaji and Majeed, 2024).		2384

2385	Human Evaluation. Human review verifies algorithmic rationality and Logical Coherence .	E Extended Discussion on Future	2432
2386	Experts audit code for theoretical soundness, examining subtle edge cases that automated suites might miss (Dunivin, 2024; Li et al., 2025c, 2022; Austin et al., 2021).	Directions	2433
2387		This appendix provides a detailed analysis of the four strategic frontiers for CoT evaluation outlined in Section 6.	2434
2388			2435
2389			2436
2390			
2391	D.4 Dialogue & Multi-turn Reasoning	E.1 Hybrid Evaluation Paradigms and Standardization	2437
2392	D.4.1 Benchmarks		2438
2393	In multi-turn dialogue, Consistency is dominant. MT-Bench (Zheng et al., 2023) and MINT (Wang et al., 2024f) assess stability. MuTual (Cui et al., 2020) and DREAM (Sun et al., 2019) focus on logic. QReCC (Anantha et al., 2021) and CICEO (Ghosal et al., 2022) target state tracking and causal inference, while MUSR (Sprague et al., 2024) and ClemBench (Chalamalasetti et al., 2023) test long-context rules.	The dichotomy between expensive human annotation and potentially biased automated evaluation remains a bottleneck. Future work should pivot towards hybrid human-AI paradigms, where models perform initial filtering or semi-automated scoring to reduce workload, followed by targeted human verification of uncertain cases (Diao et al., 2024; Li et al., 2023a). This approach mitigates the high cost and subjectivity of manual review while maintaining reliability. Furthermore, to address the lack of standardization across studies, the community must establish unified annotation guidelines and consistent evaluation protocols. Research into adversarial robustness is also vital, employing techniques to generate adversarial samples that stress-test the stability of reasoning evaluation metrics against subtle prompt variations (Cui et al., 2025; Sclar et al., 2024; Zhou et al., 2024).	2439
2394			2440
2395			2441
2396			2442
2397			2443
2398			2444
2399			2445
2400			2446
2401			2447
2402	D.4.2 Evaluation Methodologies		2448
2403	Automated Evaluation. Protocols decompose structure to verify Logical Coherence . The Q^2 framework (Honovich et al., 2021) generates constituent questions for factual support. DNLI (Welleck et al., 2019) automates Consistency checks by classifying logical relationships between utterances.		2449
2404			2450
2405			2451
2406			2452
2407			2453
2408			2454
2409			2455
2410			2456
2411	Human Evaluation. Human evaluation audits cross-turn Consistency . Experts ensure coherent belief states without violating constraints. Zheng et al. (2023) formalize this in MT-Bench using pairwise comparisons to penalize instruction drift (Suhr et al., 2021).	E.2 Multimodal and Multilingual Frontiers	2457
2412			2458
2413			2459
2414			2460
2415			2461
2416	D.5 Scientific Reasoning	As models expand beyond text, evaluation frameworks must adapt to Multimodal CoT . This involves developing MLLM-as-a-Judge frameworks capable of assessing reasoning that intertwines visual perception with textual logic, ensuring reliability under distribution shifts (Jiang et al., 2025a; Chen et al., 2024b; Wu et al., 2025a). Simultaneously, the Anglocentric bias in current benchmarks necessitates a shift towards multilingual evaluation, particularly for low-resource and minority languages. Initiatives like the Aya project and multilingual benchmarks (Bang et al., 2023; Athiwaratkun et al., 2023) demonstrate the importance of instruction-tuned evaluation across diverse linguistic landscapes to ensure equitable reasoning capabilities.	2462
2417	Scientific evaluation places a premium on Correctness and Logical Coherence . GPQA (Rein et al., 2024) sets a PhD-level standard, while MegaScience (Fan et al., 2025) addresses data scarcity. SCIBench (Wang et al., 2024e) and OlympicArena (Huang et al., 2024) test cognitive capabilities. Specialized benchmarks include ChemVTS-Bench (Huang et al., 2025b) and ChemBench (Walker et al., 2010) for chemistry, LAB-Bench (Laurent et al., 2024) for biology, and SciCode (Tian et al., 2024) for physics. Multimodal frontiers like ScienceQA (Saikh et al., 2022) and SciEval (Sun et al., 2024) lead to DiscoveryBench (Majumder et al., 2024), challenging agents to formulate hypotheses.		2463
2418			2464
2419			2465
2420			2466
2421			2467
2422			2468
2423			2469
2424			2470
2425			2471
2426			2472
2427			2473
2428		E.3 Quantifying Interpretability and Trustworthiness	2474
2429			2475
2430		Moving beyond vague notions of explainability, future research must mathematically quantify the structural validity of reasoning. Promising directions include leveraging Causal Graphs and Knowl-	2476
2431			2477
			2478
			2479

2480 edge Graphs to map the logical flow of CoT, trans-
2481 forming text-based traces into verifiable structural
2482 paths (Amayuelas et al., 2025; Nguyen et al., 2024;
2483 Zhao et al., 2025b; Luo et al., 2024). Addition-
2484 ally, integrating Explainable AI techniques can fa-
2485 cilitate human-AI collaborative review, allowing
2486 evaluators to pinpoint exactly which part of the
2487 input context influenced a specific reasoning step,
2488 thereby addressing the challenge of disguised ac-
2489 curacy where models arrive at correct answers via
2490 flawed logic (Thakkar et al., 2023; Wang and Xu,
2491 2025; Turpin et al., 2023).

2492 CoT evaluation should also transition from gen-
2493 eral benchmarks to high-stakes vertical applica-
2494 tions. In domains such as healthcare, law, and edu-
2495 cation, evaluation metrics must prioritize trustwor-
2496 thiness, evidence retrieval, and safety over simple
2497 accuracy. For instance, legal reasoning requires val-
2498 idating adherence to statutory interpretations (Yu
2499 et al., 2025), while educational applications de-
2500 mand iterative self-correction and feedback mech-
2501 anisms for pedagogical effectiveness (Jiang et al.,
2502 2024b; Wang et al., 2024d). Moreover, bias anal-
2503 ysis is crucial; frameworks must be developed to
2504 detect deep-seated stereotypes (e.g., gender or cul-
2505 tural bias) within reasoning chains (Shaikh et al.,
2506 2023; Wu et al., 2024b). In enterprise contexts,
2507 evaluation protocols must also align with auditabil-
2508 ity and compliance standards to support safe large-
2509 scale deployment.

2510 **E.4 Probing Reasoning Upper Bounds**

2511 Current benchmarks largely measure static in-
2512 ference, overlooking the potential of test-time
2513 scaling (TTS) where models simulate deliberate
2514 thinking (Snell et al., 2024; Goyal et al., 2023).
2515 Future evaluation must quantify the efficiency-
2516 performance trade-off, verifying if extended rea-
2517 soning time and search depth yield proportional
2518 accuracy gains (Brown et al., 2024; Wu et al.,
2519 2025b). This shift necessitates new benchmarks
2520 like JETTS to ensure automated judges can reli-
2521 ably monitor these complex processes (Lightman
2522 et al., 2024; Uesato et al., 2022). Simultaneously,
2523 protocols must expand to multi-agent systems, as-
2524 sessing whether increasing ensemble size follows
2525 predictable scaling laws or encounters diminishing
2526 returns due to redundancy (Li et al., 2024a; Zhuge
2527 et al., 2023).

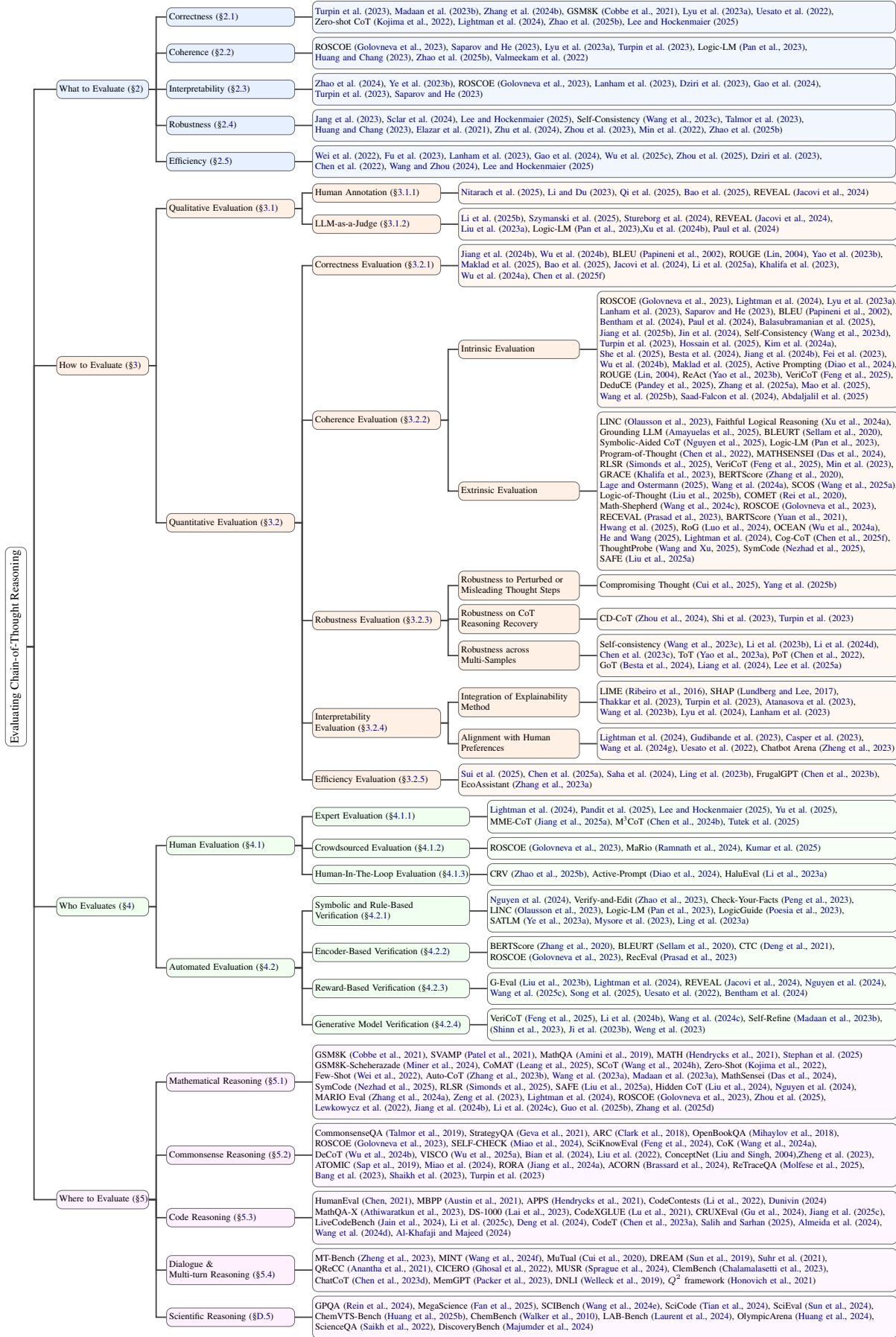


Figure 6: Our full taxonomy for the survey on Evaluating Chain-of-Thought Reasoning.

Table 1: Overview of benchmarks used in CoT evaluation across different tasks.

Task	Dataset	Description
Mathematical Reasoning	GSM8K (Cobbe et al., 2021)	8.5K grade school math word problems with step-by-step solutions.
	SVAMP (Patel et al., 2021)	Challenge set testing robustness to superficial keywords.
	MathQA (Amini et al., 2019)	GRE-level math problems with annotated operation programs.
	MATH (Hendrycks et al., 2021)	12,500 competition math problems (AMC10/12, AIME).
	AIME (Veeraboina, 2023)	High-difficulty competition problems (often included/covered by MATH).
	GSM8K-Scheherazade (Miner et al., 2024)	Synthetic extension of GSM8K to counter saturation.
Commonsense Reasoning	CommonsenseQA (Talmor et al., 2019)	Multiple-choice questions requiring commonsense knowledge.
	StrategyQA (Geva et al., 2021)	Implicit reasoning steps for multi-hop reasoning.
	ARC (Clark et al., 2018)	Science exam questions testing genuine understanding.
	OpenBookQA (Mihaylov et al., 2018)	Questions requiring knowledge retrieval and reasoning.
	ROSCOE (Golovneva et al., 2023)	Metrics for semantic consistency in reasoning chains.
	SELF-CHECK (Miao et al., 2024)	Automated step-wise error detection framework.
	SciKnowEval (Feng et al., 2024)	Scientific knowledge evaluation with knowledge triples.
	CoK (Wang et al., 2024a)	Chain-of-Knowledge framework for evidence retrieval.
	DeCoT (Wu et al., 2024b)	Causal intervention framework for bias mitigation.
	VISCO (Wu et al., 2025a)	Multimodal visual reasoning evaluation.
Code Reasoning	HumanEval (Chen, 2021)	164 hand-written programming problems with pass@k metrics.
	MBPP (Austin et al., 2021)	Basic Python problems dataset.
	APPS (Hendrycks et al., 2021)	Competition-level programming problems.
	CodeContests (Li et al., 2022)	Programming competition problems from various platforms.
	MBXP (Athiwaratkun et al., 2023)	Multilingual programming benchmarks.
	MathQA-X (Athiwaratkun et al., 2023)	Multilingual MathQA for code reasoning.
	DS-1000 (Lai et al., 2023)	Data science programming benchmark.
	CodeXGLUE (Lu et al., 2021)	Code understanding and generation benchmark.
	CRUXEval (Gu et al., 2024)	Execution reasoning evaluation.
LiveCodeBench (Jain et al., 2024)	Contamination-free code evaluation.	
Dialogue & Multi-turn	MT-Bench (Zheng et al., 2023)	Multi-turn dialogue evaluation.
	MINT (Wang et al., 2024f)	Multi-turn instruction following.
	MuTual (Cui et al., 2020)	Multi-turn dialogue reasoning.
	DREAM (Sun et al., 2019)	Dialogue-based reading comprehension.
	QReCC (Anantha et al., 2021)	Question rewriting for conversational context.
	CICERO (Ghosal et al., 2022)	Conversational inference and causal reasoning.
	MUSR (Sprague et al., 2024)	Multi-turn soft reasoning.
	ClemBench (Chalamalasetti et al., 2023)	Game-theoretic rule adherence.