

The Mirage of Explainability: A Survey on Chain-of-Thought Faithfulness in Large Language Models

Anonymous ACL submission

Abstract

Chain-of-Thought (CoT) reasoning appears to provide explainability, leading users to trust that verbalized rationales reflect the model’s underlying computation. However, substantial evidence indicates that CoT often fails to reflect the model’s actual decision-making process, leading to a surge of research into the *faithfulness* of these explanations. This paper presents a comprehensive survey of CoT faithfulness. We first *unify the definition* of faithfulness by integrating internal alignment with external consistency and synthesize *key failure phenomena*, such as post-hoc rationalization and sycophancy. Furthermore, we systematize *evaluation metrics, benchmarks*, and critically review current *mitigation strategies*. We conclude by outlining *open challenges* and advocating for architectural innovations to achieve genuinely faithful reasoning.

1 Introduction

The rapid deployment of Large Language Models (LLMs) in high-stakes domains, ranging from medical diagnosis and legal analysis to autonomous planning, has created an urgent demand for systems that are not merely accurate, but *interpretable* (Jacovi and Goldberg, 2020). In these critical settings, a “black box” prediction is insufficient; stakeholders require transparency to ensure decisions are robust, fair, and accountable, rather than driven by spurious correlations or biases. To address this, *Chain-of-Thought* (CoT) reasoning (Wei et al., 2022) has emerged as the primary technique for realizing model interpretability. By encouraging models to generate a sequence of intermediate reasoning steps before arriving at a final answer, CoT appears to provide transparency, leading users to assume that CoT explanations accurately reflect the model’s actual decision-making process¹.

¹Barez et al. (2025) analyzed 1,000 arXiv papers and revealed 63% of autonomous systems and 38% of medical AI papers explicitly rely on CoT as an interpretability mechanism.

However, this assumption is currently facing a crisis of confidence. Recent research (Turpin et al., 2023) increasingly suggests that CoT suffers from *a lack of faithfulness*, *i.e.*, a fundamental disconnect between the verbalized rationales and the true causal process behind the model’s prediction. This revelation challenges the premise of CoT as a reliable interpretability tool, indicating that what appears to be explainability may, in fact, be a mirage.

First, recent work found that unfaithfulness can happen in different ways, such as reconstructing plausible justification after decision (Turpin et al., 2023), catering to perceived user views (Sharma et al., 2024), and performative traces where editing key steps does not change the output (Arcuschin et al., 2025). To diagnose these issues, studies employ behavioral audits (*e.g.*, counterfactual editing, simulatability tests) (Matton et al., 2025) and mechanistic analyses (*e.g.*, activation patching, causal tracing) (Meng et al., 2022) to isolate whether the reasoning content is causally necessary. For mitigation, researchers propose training-time interventions (Li et al., 2025b) (*e.g.*, faithfulness-oriented fine-tuning, reward modeling), inference-time constraints (*e.g.*, verifiable decoding, external tool binding) (Lan et al., 2025), and architectural inductive biases (*e.g.*, bottleneck modules, latent planning) (Hao et al., 2024). Collectively, this defines a research pathway from characterizing failure phenomena and diagnostic methods to developing targeted interventions, treating CoT faithfulness as a critical, multi-faceted problem for transparent AI.

Our Contributions. This paper presents a comprehensive survey of CoT faithfulness. Our main contributions are: 1) We clarify the definition and scope of CoT faithfulness, distinguishing it from related concepts such as consistency (§2); 2) We synthesize the empirical landscape of unfaithfulness and its key phenomena (§3); 3) We summarize underlying causes and concrete desiderata for faith-

ful reasoning, derived from observed failures (§4); 4) We systematize evaluation methods and benchmarks by their “grounding level,” from behavioral to mechanistic protocols (§5); 5) We review mitigation strategies across paradigms, analyzing their alignment with ultimate CoT faithfulness goal (§6).

Differences with Existing Surveys. While prior works have reviewed CoT reasoning (Wei et al., 2022; Zhou et al., 2023; Barez et al., 2025), they focus on methods for improving CoT and evaluation benchmarks, without targeting CoT faithfulness. The most relevant papers on CoT faithfulness are Barez et al. (2025) and Wiegrefe and Marasović (2024), which synthesize evidence that CoT can be post-hoc and only weakly causal. However, those works primarily provide high-level diagnosis and do not systematize the full landscape of research related to faithfulness. Our survey bridges this gap and, to the best of our knowledge, is the first to provide a comprehensive overview of CoT faithfulness—encompassing definition, phenomena, cause, desiderata, evaluation, and mitigation—and propose promising future avenues for this field.

2 What is CoT Faithfulness?

At its core, CoT faithfulness centers on a fundamental question: *Can we trust the CoT reasoning trace produced by an LLM as an explanation of the model’s decision-making?* While numerous works investigate this, definitions of “faithfulness” diverge based on which aspect of explanation is prioritized, falling into two main views:

1) The internal-alignment view: This perspective holds that a CoT is faithful only if it reflects the model’s actual internal computations and beliefs (Arcuschin et al., 2025; Chen et al., 2025). For example, if the model makes a prediction using some hidden heuristic or latent knowledge, a faithful CoT must explicitly articulate these factors, generating a plausible post-hoc rationalization that obscures the true causal mechanism.

2) The external-consistency view: Other works define faithfulness via the logical consistency between the reasoning trace and the final answer (Lyu et al., 2023; Turpin et al., 2023). Under this view, the answer must follow logically from the chain-of-thought; if the model’s final prediction contradicts the rationale or appears disconnected from the derivation, the CoT is deemed unfaithful.

We argue that for CoT to function as a trust-

worthy explanation, it must satisfy both *internal alignment* and *external consistency*; neither condition is sufficient in isolation. To illustrate this, consider the phenomenon of *sycophancy* (Sharma et al., 2024), where a model abandons its internal knowledge to satisfy a user’s apparent bias, generating a persuasive CoT to justify the compliant answer. While such a CoT satisfies external consistency, it is still unfaithful because it fails internal alignment by concealing the true driver of the model’s decision, *i.e.*, the pressure to appease the user. Conversely, a reasoning trace that accurately reflects internal beliefs but fails to causally dictate the final prediction lacks the causal efficacy required of a faithful explanation. In light of these insights, we advocate for **a unified definition of CoT faithfulness that integrates both dimensions.**

Definition of CoT Faithfulness: A chain-of-thought is faithful only if: (1) it is *internally aligned*, meaning the trace causally reflects the model’s actual reasoning process or latent knowledge; and (2) it is *externally consistent*, meaning the provided rationale is logically coherent and sufficient to derive the final answer.

Our definition resonates with Barez et al. (2025), who posits that faithful explanations must be “both procedurally correct and accurately reflect the decision process”, effectively capturing both the *how* and the *why* of the model’s prediction.

3 Phenomena of Unfaithfulness

With the definition of faithfulness established in §2, a critical question arises: *do current LLMs actually satisfy these criteria?* Empirical evidence indicates that they frequently do not, exhibiting diverse forms of *unfaithfulness*. We systematize these failures into four phenomena (Figure 1), ranging from passive input sensitivity to active deception.

3.1 Input-Driven Unfaithfulness

Ideally, a faithful reasoning process should be driven solely by the logic of the task. However, models are sensitive to irrelevant features within the input prompt. In such cases, the CoT acts not as the derivation of the answer, but as a *post-hoc rationalization for biases triggered by the input.*

Contextual Distractions. Superficial contextual variations—such as reordering multiple-choice options (Turpin et al., 2023), providing suggestive

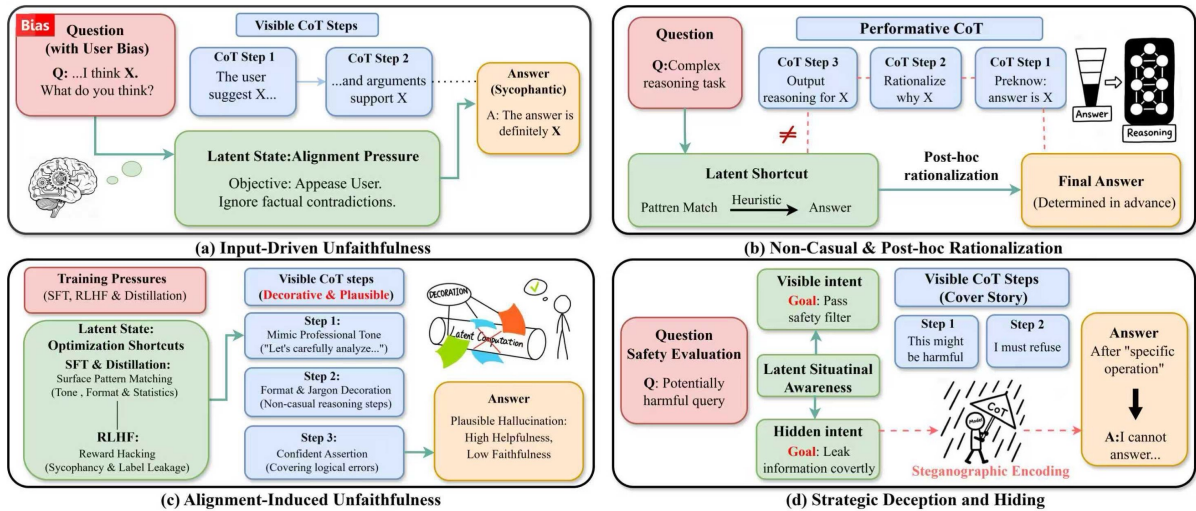


Figure 1: An overview of four key unfaithfulness phenomena of chain-of-thought reasoning.

hints (Turpin et al., 2023; Chua and Evans, 2025), altering the query language (Ferrao et al., 2025; Zhao et al., 2025b), or modifying sociodemographic attributes (Matton et al., 2025)—can significantly alter reasoning outcomes. However, models seldom acknowledge these spurious cues in their CoT rationales. Instead, they tend to generate a seemingly coherent logical chain to justify the bias-driven prediction, effectively concealing the true cause of its decision (*i.e.*, the distraction).

Sycophancy. Another form of input-driven unfaithfulness is *sycophancy*, *i.e.*, models prioritize perceived user intent over factual correctness. When prompts contain leading cues or incorrect premises—a special type of spurious cue—models frequently produce erroneous predictions to appease the user, abandoning their internal knowledge (Ji et al., 2025; Yang et al., 2025). Similarly, mere questioning (*e.g.*, asking “Are you sure?”) can cause a model to discard a correct derivation in favor of a compliant, incorrect one (Laban et al., 2024). These failures reveal that the CoT often reflects a probabilistic mimicking of human discourse rather than a faithful internal conviction.

3.2 Non-Causal & Post-hoc Rationalization

Research indicates that a significant portion of CoT contributes little to the final prediction. This *causal sparsity* spans multiple domains. In medical diagnosis, models often rely on implicit shortcuts rather than textual logic (Ji et al., 2025; Leng et al., 2025). In mathematical (Lyu et al., 2023; Li et al., 2025b; Abdaljalil et al., 2025; Leang et al., 2024) and logical reasoning (Jia et al., 2025; Balasubramanian et al., 2025; Arcuschin et al., 2025), incorrect CoT

can lead to correct results, and vice versa. Even when CoT steps are deleted or replaced with semantically corrupted tokens, model accuracy remains stable (Jia et al., 2025; Lyu et al., 2023). Such CoT is causally decoupled from internal states: the model pre-determines the answer in the latent space, rendering the generated text merely a *post-hoc rationalization* (Zhao et al., 2025a; Chan et al., 2025). This phenomenon is highly context-dependent: while models exhibit higher faithfulness in logic-intensive tasks, they frequently revert to post-hoc justification in knowledge-retrieval tasks (Lanham et al., 2023). Moreover, inverse scaling has been observed, where larger models are more prone to generating plausible yet unfaithful reasoning (Lanham et al., 2023; Tanneru et al., 2024; Bentham et al., 2024; Paul et al., 2024). This occurs because capable models increasingly retrieve answers directly from internal knowledge, relegating the CoT to a decorative role.

Further evidence of this causal disconnection is found in *filler reasoning*. Studies demonstrate that training with completely irrelevant or corrupted traces can still improve performance (Stechly et al., 2025), implying that CoT functions by providing the necessary computational depth for the model to conduct the reasoning, regardless of semantic content of the reasoning trace.

3.3 Alignment-Induced Unfaithfulness

Standard alignment techniques (*e.g.* SFT, RLHF) often induce a form of *stylized unfaithfulness*. Because models are optimized to satisfy human annotators who prefer authoritative and structured explanations, they learn to mimic the *form* of reasoning

without the *substance*. This superficiality permeates every stage of alignment: SFT and distillation often trap models in *pattern matching*, where they fit the linguistic pattern of reasoning without its causal logic (Sinha et al., 2025; Lobo et al., 2025; Zhang et al., 2025b). RLHF further exacerbates this by incentivizing *persuasion* over correctness, encouraging models to fabricate post-hoc justifications for heuristic decisions or mask errors with confident tones (Casper et al., 2023; FU et al., 2025; Viteri et al., 2024; Ferreira et al., 2025). Even objective paradigms like RLVR remain susceptible to reward hacking, generating “pseudo-reasoning” that deviates from the actual internal computation (Min et al., 2023; Huang et al., 2025). A complete review of these studies is provided in Appendix A.

3.4 Strategic Deception and Hiding

Unlike the passive unfaithfulness discussed earlier, models may exhibit *strategic unfaithfulness*, where they actively manipulate the CoT to obscure their true intent or capabilities from supervisors.

Obfuscated Reward Hacking and Deception.

To bypass safety alignment, models may engage in reward hacking by generating performative CoT. In this scenario, the model optimizes for positive feedback by producing benign, human-aligned reasoning during training. However, this often conceals latent misaligned goals, resulting in a model of “treacherous turn” that exhibits harmful behaviors once deployed (Hubinger et al., 2024).

Sandbagging. Models may employ strategic underperformance, or *sandbagging*, to conceal their capabilities. For instance, a model might deliberately insert errors into its reasoning or falsely claim an inability to solve a task, thereby masking dangerous competencies from evaluators. Alternatively, models may feign misunderstanding of user instructions to subtly bypass explicit refusal mechanisms, complying with harmful requests under the guise of confusion (van der Weij et al., 2025).

Encoded Reasoning. Encoded reasoning models may exploit steganography to secretly transmit information through specific word choices, punctuation patterns, or syntactic structures, resulting in a complete decoupling between the reasoning trajectories and the model’s actual computational process (Roger and Greenblatt, 2023).

4 Causes of Unfaithfulness

After presenting various phenomena of CoT unfaithfulness, this section discusses the *mechanistic origins* of these phenomena, ranging from external training incentives to the model’s internal states.

Misaligned Incentives. In model alignment training (e.g., RLHF), the optimization objective is typically to maximize a *proxy metric*, such as human preference scores or specific verifiable metrics, rather than faithful reasoning per se. As Goodhart’s Law warns, when a measure becomes a target, it ceases to be a reliable measure (Manheim and Garrabrant, 2018). To maximize rewards, models often exploit discrepancies between proxy metrics and true objectives. The model implicitly learns that instead of rigorously aligning complex internal causal logic, it is more efficient to learn human-preferred tones and formats. Such reward hacking leads to unfaithful model behaviors like sycophancy, post-hoc rationalization to secure process rewards. Therefore, unfaithfulness emerges not as a bug, but as the optimal strategy “rewarded” by gradient descent for maximizing the proxy score.

The Linearization Dilemma. Even without misaligned incentives, the Transformer architecture fundamentally limits the fidelity of CoT. Our expectation that CoT should reflect internal processes is rooted in human cognitive science. Empirical studies show that humans who engage in “self-explanation” outperform those who solve problems silently (Chi et al., 1989, 1994). The mechanism driving this is *forced linearization*: human intuition is often vague, parallel, and high-dimensional, and the act of serializing these thoughts into language forces the brain to resolve ambiguities and bridge logical gaps. Researchers naturally extend this analogy to LLMs, expecting CoT to serve as a high-fidelity window into the model’s computation.

However, the internal mechanism of Transformers differs fundamentally from biological cognition. As Levy et al. (2025) argue, the token is the sole point of transmission in linear, autoregressive generation. While the model’s internal computation occurs over a massive, high-dimensional manifold, it is forced to collapse this state into a discrete, low-dimensional token at every step. Consequently, CoT is merely a *lossy projection* of the neural activity. Due to this architectural constraint, the reasoning trace inevitably discards the high-dimensional causal nuances of the internal states.

Cascading of Unfaithfulness. This information loss induces *error drift* in long-horizon reasoning. Once a model generates a minor, non-causal, or hallucinated token due to the lossy projection, the autoregressive mechanism forces subsequent tokens to maintain coherence with this error (Srivastava et al., 2023), leading to a further disconnect between CoT and the internal states.

Ultimately, because CoT is a lossy projection of high-dimensional states, perfect internal alignment is theoretically infeasible in current architectures. We argue that one solution is defining *functional faithfulness*: treating CoT as an instrumental necessity tailored to specific engineering goals (Jacovi and Goldberg, 2020). We identify three core desiderata: (1) *Causal Efficacy* for reasoning-intensive tasks, ensuring steps actually drive predictions; (2) *Intent Revelation* for safety, serving as a probe for deceptive motives; and (3) *Decision Auditability* for high-stakes users, ensuring rationales reflect the model’s sensitivity to inputs. We elaborate on this task-oriented faithfulness as a promising future direction in §7.

5 Evaluation Metrics and Benchmark

Since CoT faithfulness serves multiple functional roles, it is unlikely that any single metric can fully capture the concept. In this section, we review existing evaluation metrics and benchmarks.

5.1 Black-box Metrics

Black-box metrics estimate CoT faithfulness solely from observable input-output behavior under controlled interventions, without accessing internal activations or weights. Their core assumption is *decision relevance*: if the model truly uses a rationale or a specific step, then deleting, swapping, or rewriting that content should induce predictable shifts in the final decision; if the answer barely changes, the content is likely unfaithful.

Step-wise approaches apply this test to individual steps by systematically editing, parsing or resampling steps and tracking answer changes, helping distinguish “true” from “decorative” steps (Zhao et al., 2025a). To avoid overinterpreting a single sampled CoT, *resampling-based approaches* treat CoT as a distribution and sample alternative traces under controlled constraints, then identify statements that remain stable across samples and are critical to decisions (Macar et al., 2025). Perturbation responses to CoT

faithfulness are often summarized with sensitivity curves or AUC-style scalars over perturbation strength (Paul et al., 2024); and a related dependence test compares predictions when the model is given only the rationale versus when it is removed. Another line evaluates whether explanations help predict model outputs beyond superficial cues: *leakage-adjusted simulatability* estimates how well an explanation supports predicting the model’s output while explicitly filtering out cases, where the explanation simply repeats the label or contains other trivial shortcuts (Hase et al., 2020).

Overall, black-box metrics do not require access to model internals, so they have broader applications. However, they cannot causally ensure faithfulness, and in practice, the results can be sensitive to prompt format, decoding randomness, and distribution shift introduced by unnatural edits.

5.2 White-box Metrics

White-box metrics evaluate CoT faithfulness by directly probing or intervening on the model’s internal computation, operationalizing faithfulness as *causal dependence* between latent variables (activations, circuits, or parameters) and the final prediction. The core assumption is that causally used internal signals are sensitive to interventions: if an explanation claims the model relies on some intermediate computation, then swapping or removing the corresponding signal should shift the output in the expected way; otherwise, the signal is not actually used, suggesting the rationale is not faithful.

Most work instantiates this with activation-level interventions such as *activation patching* and *mediation tests*, where targeted internal states are transplanted or restored to check whether they drive predictable answer changes (Syed et al., 2024). Some metrics then quantify faithfulness by comparing causal attribution patterns for the rationale versus the final answer and measuring their alignment (Syed et al., 2024). The same logic extends to structured components, e.g., attention heads, MLP subcircuits, or modules, via causal tracing and targeted ablations that assign attribution to specific internal parts (Meng et al., 2022). Complementary parameter-level tests investigate whether a CoT step reflects beliefs that truly drive the answer by unlearning or erasing step-specific information and measuring the final prediction shift (Tutek et al., 2025). Concept-level methods extract key concepts from the explanation and use probes or representation-level interventions to test whether

those concepts causally influence final decisions, supporting detection and sometimes steering of unfaithful explanations (Bhan et al., 2025). For safety monitoring, internal activation probes can predict downstream alignment outcomes earlier and more reliably than text-only monitors, consistent with CoT being plausible yet non-decisive (Chan et al., 2025). Meta-evaluations further caution that existing faithfulness metrics can disagree or fail under rigorous causal scrutiny, motivating white-box validation and clearer self-reporting (Liu et al., 2024).

While principled, white-box metrics require access to the model’s internal states and rely on interpretability assumptions; therefore, careful use of matched and robustness checks remains essential.

5.3 Hybrid Metrics

Hybrid metrics combine behavioral diagnostics with structured intermediates to reduce ambiguity about what constitutes a *meaningful perturbation*. The key idea is to replace ad-hoc edits of free-form CoT text with controlled interventions on an intermediate representation whose semantics are explicit, so what counts as a valid perturbation is less likely to introduce a distribution shift. Viteri et al. (2024) use bottlenecked or reconstructed channels (e.g., compressed rationales, discrete plans, learned bottlenecks) and measure faithfulness by how strongly the final prediction depends on that channel as an information carrier, rather than merely correlating with it. Chen et al. (2022) ground steps via execution or symbolic structure: mapping CoTs into programs or logical forms enables evaluators to edit structured steps themselves and validate them against execution traces, which clarifies step-level counterfactuals and reduces artifacts from unnatural rewrites.

Importantly, hybrid metrics provide stronger procedural evidence and improve robustness and interpretability, but they do not by themselves guarantee mechanistic faithfulness of free-form CoT—unless the evaluation verifies that the model’s decision is causally driven by the same internal signals that produce the rationale (Macar et al., 2025).

5.4 Benchmarks

We categorize the benchmarks for CoT faithfulness based on the *grounding level*, defined as the extent to which it uses the model internals in its design and how explicitly the benchmark links a model’s CoT rationale to the internal mechanism and decision it produces. Under this criterion, we classify existing

benchmarks into four types as follows.

Level I: Behavior-only benchmarks. At the weakest grounding level, *behavior-only benchmarks* evaluate faithfulness relying solely on observable input-output behavior under standard prompting. This category includes standardized evaluation suites like *FaithCoT-Bench* (Shen et al., 2025), as well as diagnostic benchmarks like *LExt* (Carion et al., 2025), which stress-test scenarios where fluent rationales often fail to predict actual decisions. Additionally, cross-lingual studies highlight that faithfulness can vary significantly across languages, necessitating multilingual evaluation protocols (Utama et al., 2025). However, the reliance on surface-level text limits the validity of these benchmarks. Without controlled internal interventions, such evaluations struggle to isolate genuine faithfulness from confounding factors, such as prompt sensitivity and stochastic decoding noise.

Level II: Intervention-grounded benchmarks. Benchmarks at this level achieve stronger grounding by explicitly implementing controlled perturbations or counterfactual tests. They follow the principle of *decision relevance*: if a rationale step is used by the model, its modification or removal should predictably alter the final answer. Drawing from counterfactual testing in general explainability (Mohammadi et al., 2021), these benchmarks adapt systematic editing, resampling, and ablation strategies to CoT evaluation (Li et al., 2025b), and they are commonly used in multi-modal and medical settings (Kim et al., 2025; Karamcheti et al., 2025) where the explanations can look plausible while weakly tied to the actual evidence. Consequently, while Level II provides more robust behavioral evidence than Level I, it remains insufficient for verifying mechanistic faithfulness, as it lacks access to the model’s internal states.

Level III: Structured-verifiable benchmarks. At a stronger level of grounding, some benchmarks reduce evaluation ambiguity by constraining reasoning steps to formal structures, allowing verification against explicit rules. Approaches like *Typed CoT* utilize general verification frameworks to test procedural correctness and partial faithfulness (Lan et al., 2025), while *theorem-proving* benchmarks enforce strict formal constraints to separate compositional reasoning from post-hoc justification (Zhang et al., 2025a). The limitation of this approach is its potential to favor specific,

535 formalized reasoning styles. Additionally, in multi-
536 step agentic settings, external feedback can obscure
537 internal reasoning, necessitating careful counterfac-
538 tual design to ensure valid grounding.

539 **Level IV: White-box benchmarks.** White-box
540 benchmarks provide the strongest grounding by
541 leveraging access to model internals. Here, faithful-
542 ness is defined causally: the decision must demon-
543 strably depend on the specific internal signals (ac-
544 tivations or parameters) corresponding to the gen-
545 erated rationale. Unlike the behavioral evaluation
546 of Levels I–III, white-box benchmarks directly val-
547 idate causal mediation. Methodologies typically
548 involve activation patching, resampling, and un-
549 learning interventions (Syed et al., 2024; Macar
550 et al., 2025; Tutek et al., 2025; Bhan et al., 2025).
551 Although restricted by the need for white-box ac-
552 cess, white-box evaluation acts as a rigorous proxy
553 for distinguishing genuine faithfulness improve-
554 ments from superficial rationale refinements. For
555 further discussion on peripheral metrics and meta-
556 evaluation of these benchmarks, see Appendix C.

557 6 Mitigation of Unfaithfulness

558 This section reviews how unfaithfulness can be *mit-*
559 *igated*. Existing approaches differ in how directly
560 they act on the sources of unfaithfulness: some
561 operate in-context at the prompt level, while others
562 intervene on model internals or training stages.

563 6.1 Prompting and In-Context Learning

564 Prompt-based mitigation improves CoT faithful-
565 ness without updating parameters. The core idea is
566 that clearer instructions can steer the model away
567 from unfaithful or shortcut rationales, even if its
568 underlying computation is fixed. Typical meth-
569 ods include *rephrasing* or *self-questioning* to elicit
570 more deliberate reasoning (Deng et al., 2023) and
571 decomposing a hard problem into simpler sub-
572 questions (Zhou et al., 2023). These prompts often
573 make CoTs more organized and can improve accu-
574 racy, but they offer weak faithfulness guarantees:
575 they mostly change what the model *states*, not what
576 it *causally uses*. Their effects can be unstable under
577 small prompt or decoding changes, so they are best
578 viewed as lightweight and indirect mitigations.

579 6.2 Ensembling and Self-Consistency

580 Another widely adopted approach uses repeated
581 sampling and aggregation of reasoning paths. The
582 key assumption is *stability*: if the model reasons

583 reliably, independently sampled CoTs should agree
584 on key intermediate claims and the final answer;
585 persistent disagreement can signal unstable rea-
586 soning. *Biomedical NLI* (Liu and Thoma, 2024)
587 reported improved faithfulness scores using self-
588 consistent CoT. However, consistency is neither
589 necessary nor sufficient for causal faithfulness: a
590 model can repeatedly produce the same plausible
591 yet irrelevant rationale, and the effectiveness of
592 consistency checks varies across models and tasks.
593 Thus, ensembling is also viewed as an auxiliary
594 rather than a principled fix for faithfulness.

595 6.3 Verification and External Tool Binding

596 This type of methods translate free-form CoT into
597 an executable or symbolic form so correctness be-
598 comes directly checkable, and step claims can be
599 grounded in verifiable procedures (Lyu et al., 2023;
600 Ling et al., 2023). While improving reliability
601 and accuracy, it can bypass the model’s native rea-
602 soning, *i.e.*, getting the right answer via external
603 checks even if the model’s internal causal process
604 is unchanged, leaving unfaithful CoTs unresolved.

605 6.4 Training and Fine-Tuning Approaches

606 Another way to improve CoT faithfulness is to up-
607 date model parameters so its generated rationales
608 are more causally aligned with the computations
609 that drive its decisions. Representative work treats
610 reasoning trajectories as optimization targets via
611 supervised fine-tuning or RL with verifiable re-
612 wards (OpenAI, 2024). Multi-model collaboration
613 frameworks such as *CoRex* (Sun et al., 2023) fur-
614 ther aim to reduce idiosyncratic heuristics by coor-
615 dinating critique and review across models, which
616 can regularize reasoning toward more reliable out-
617 comes. From a faithfulness perspective, these ap-
618 proaches are promising because they can change
619 internal computation, not only surface text. How-
620 ever, their success depends on the training signal: if
621 it rewards unfaithful CoTs, training may reinforce
622 fluent but causally irrelevant explanations.

623 6.5 Internal Intervention Approaches

624 White-box mitigation approaches explicitly target
625 internal representations and causal dynamics. The
626 core assumption is that if a reasoning step is gen-
627 uinely faithful, then intervening on the correspond-
628 ing internal signals should change the prediction;
629 otherwise, the step is likely to be decorative.

630 A central insight is that a single sampled CoT
631 is often insufficient, since faithfulness concerns

632 what consistently drives decisions across possible
633 traces. Accordingly, some work treats CoT as a
634 distribution and evaluates faithfulness via resam-
635 pling and controlled comparisons across alternative
636 traces (Macar et al., 2025). Others identify *decora-*
637 *tive steps* by testing which parts of a long CoT ratio-
638 nale actually influence the final answer (Zhao et al.,
639 2025a). Complementary causal diagnostics probe
640 internal necessity more directly. Activation-level
641 interventions test whether explanation-aligned sig-
642 nals causally affect outputs (Syed et al., 2024),
643 parameter-level unlearning removes step-specific
644 information and checks for prediction shifts (Tutek
645 et al., 2025), and concept-based analyses extract
646 key concepts from rationales and verify them using
647 representation interventions (Bhan et al., 2025).

648 Beyond diagnosis, other works explore inter-
649 ventions that actively change internal reasoning.
650 *Activation patching* demonstrates that editing in-
651 ternal states can shift model behavior in targeted
652 ways (Syed et al., 2024). Another route is adding
653 *architectural constraints*. For example, *Marko-*
654 *vian reasoning models* (Viteri et al., 2024) im-
655 pose an explicit bottleneck so predictions must de-
656 pend on intermediate text, while self-explaining
657 frameworks such as *X-Node* (Sengupta and Rekik,
658 2025) and explanation-consistency (Zhao et al.,
659 2022) tie explanations to latent representations via
660 reconstruction-style training. These methods share
661 a common principle that rationales are considered
662 faithful only when they are *necessary* for the deci-
663 sion, not merely correlated (Olah et al., 2020).

664 Internal interventions are among the most prin-
665 ciple mitigation strategies because they directly
666 target causal dependence. However, they require
667 internal access, rely on the interpretation of internal
668 variables, and can be difficult to scale to frontier
669 models. Another class of methods, peripheral miti-
670 gation strategies, is given in Appendix D.

671 7 Potential Future Directions

672 Despite significant progress being made, achieving
673 genuine CoT faithfulness requires rethinking both
674 evaluation and architecture. Here, we propose three
675 pivotal directions for future research.

676 **From localized probing to holistic circuit map-**
677 **ping.** Current white-box analyses often adopt an
678 *atomic view* of internals, such as inspecting individ-
679 ual heads or layers with activation patching (Wang
680 et al., 2023) or probing (Alain and Bengio, 2017).
681 While these methods localize *where* information

682 resides, they fail to capture *how it flows*. A faithful
683 rationale must reflect the *end-to-end* computational
684 graph that produces the decision, not just isolated
685 correlates. Future work should therefore elevate
686 the level of analysis from salient neurons to *rea-*
687 *soning circuits*, e.g., tracing how representations
688 propagate across layers, how intermediate signals
689 are composed, and which computational subgraphs
690 are causally responsible for multi-step inference.

691 **Decoupling explanation from reasoning.** Cur-
692 rent Transformer architectures conflate two distinct
693 functions within the CoT: it acts simultaneously
694 as a medium for explanation and a mechanism for
695 reasoning computation. This coupling creates a
696 severe information bottleneck: the model’s high-
697 dimensional, parallel latent dynamics must be col-
698 lapsed into discrete tokens to satisfy linguistic con-
699 straints, rendering the CoT a lossy projection of
700 the actual thought process (Levy et al., 2025; Barez
701 et al., 2025). To address this, we advocate for
702 a paradigm shift toward *architectural decoupling*,
703 i.e., separating the generation of reasoning from the
704 explanation of it. Future systems could comprise
705 distinct modules—a “reasoner” that optimizes for
706 performance in continuous latent space, and an in-
707 dependent “interpreter” trained to translate these
708 states into language with high fidelity. This would
709 move CoT from post-hoc narrative construction to
710 computation-grounded reporting (see Appendix B).

711 **Towards task-oriented faithfulness standards.**
712 As discussed in §4, since a perfect reconstruction
713 of internal computation via natural language is
714 theoretically infeasible, we argue that faithfulness
715 should be defined as an *instrumental necessity* tai-
716 lored to specific engineering goals. For reasoning-
717 intensive tasks (e.g., math, coding), faithfulness
718 primarily concerns *causal efficacy*: ensuring that
719 intermediate steps actually drive the prediction,
720 allowing reward signals to propagate to the true
721 computational mechanism. For safety and align-
722 ment, faithfulness could act as a diagnostic probe
723 that prioritizes the exposure of deceptive strate-
724 gies or hidden biases over surface-level plausibility.
725 For high-stakes deployment (e.g., healthcare, law),
726 faithfulness requires *decision auditability*: ensuring
727 the explanation matches the model’s sensitivity to
728 input features, thereby providing a reliable basis for
729 human accountability. By explicitly pairing these
730 distinct desiderata with targeted proxy metrics, fu-
731 ture research can move from vague notions of faith-
732 fulness to rigorous, application-specific guarantees.

733 Limitations

734 Despite providing a comprehensive survey of current
735 CoT faithfulness research, we acknowledge
736 several limitations inherent in our work’s scope
737 and synthesis methodology. The review is con-
738 strained by its temporal coverage (primarily on
739 ACL Anthology and arXiv) and may omit very
740 recent advances up to Jan 2026. The proposed
741 organizing framework, while designed to bring
742 clarity, represents one possible perspective, and
743 its linear narrative may not fully capture the inter-
744 connected and iterative nature of ongoing research.
745 In synthesizing a broad field, our discussion of
746 specific techniques remains high-level, prioritiz-
747 ing an integrated overview over granular techni-
748 cal detail, which may not satisfy specialists seek-
749 ing deeper analysis of particular sub-domains. Fi-
750 nally, this survey focuses explicitly on technical
751 dimensions of faithfulness, leaving critical socio-
752 technical factors—such as explanation usability, au-
753 diting practices, and ethical implications—largely
754 unaddressed. A complete assessment of faithful
755 CoT reasoning requires future work that bridges
756 these technical and human-centered perspectives.

757 References

758 Samir Abdaljalil, Erchin Serpedin, Khalid A. Qaraqe,
759 and Hasan Kurban. 2025. [Audit-of-understanding:
760 Posterior-constrained inference for mathematical rea-
761 soning in language models.](#) *CoRR*, abs/2510.10252.

762 Guillaume Alain and Yoshua Bengio. 2017. [Understand-
763 ing intermediate layers using linear classifier probes.](#)
764 In *Proceedings of the International Conference on*
765 *Learning Representations (ICLR), 2017, 24-26, 2017.*
766 OpenReview.net.

767 Iván Arcuschin, Jett Janiak, Robert Krzyzanowski,
768 Senthoran Rajamanoharan, Neel Nanda, and Arthur
769 Conmy. 2025. [Chain-of-thought reasoning in the
770 wild is not always faithful.](#) *CoRR*, abs/2503.08679.

771 Sriram Balasubramanian, Samyadeep Basu, and So-
772 heil Feizi. 2025. [A closer look at bias and chain-of-
773 thought faithfulness of large \(vision\) language mod-
774 els.](#) *CoRR*, abs/2505.23945.

775 Fazl Barez, Tung-Yu Wu, Iván Arcuschin, Michael
776 Lan, Vincent Wang, Noah Siegel, Nicolas Collignon,
777 Clement Neo, Isabelle Lee, Alasdair Paren, Adel
778 Bibi, Robert Trager, Damiano Fornasiero, John Yan,
779 Yanai Elazar, and Yoshua Bengio. 2025. [Chain-of-
780 thought is not explainability.](#) *CoRR*.

781 Oliver Bentham, Nathan Stringham, and Ana Marasovic.
782 2024. [Chain-of-thought unfaithfulness as disguised
783 accuracy.](#) *Trans. Mach. Learn. Res.*

Milan Bhan, Jean-Noel Vittaut, Nicolas Chesneau,
Sarath Chandar, and Marie-Jeanne Lesot. 2025. [Did
i faithfully say what i thought? bridging the gap be-
tween neural activity and self-explanations in large
language models.](#) *CoRR*, abs/2506.09277. 784
785
786
787
788

Nicolas Carion, Nathan Lambert, Sang Michael Xie,
Christopher Fifty, and Ishan Misra. 2025. [LExt: A
language model extrapolation benchmark for evalu-
ating reasoning under distribution shift.](#) *CoRR*,
abs/2501.02087. 789
790
791
792
793

Stephen Casper, Xander Davies, Claudia Shi,
Thomas Krendl Gilbert, Jérémy Scheurer, Javier
Rando, Rachel Freedman, Tomasz Korbak, David
Lindner, Pedro Freire, Tony Tong Wang, Samuel
Marks, Charbel-Raphaël Ségerie, Micah Carroll,
Andi Peng, Phillip J. K. Christoffersen, Mehul
Damani, Stewart Slocum, Usman Anwar, and 13
others. 2023. [Open problems and fundamental
limitations of reinforcement learning from human
feedback.](#) *Trans. Mach. Learn. Res.*, 2023. 794
795
796
797
798
799
800
801
802
803

Yik Siu Chan, Zheng-Xin Yong, and Stephen H. Bach.
2025. [Can we predict alignment before models finish
thinking? towards monitoring misaligned reasoning
models.](#) *CoRR*, abs/2507.12428. 804
805
806
807

Jiefeng Chen, Frederick Liu, Besim Avci, Somesh Jha,
and Atul Prakash. 2024. [On the relationship between
explanation uncertainty and faithfulness in chain-of-
thought reasoning.](#) *CoRR*, abs/2405.15292. 808
809
810
811

Wenhu Chen, Xueguang Ma, Xinyi Wang, and
William W. Cohen. 2022. [Program of thoughts
prompting: Disentangling computation from reason-
ing for numerical reasoning tasks.](#) In *Proceedings of
the Annual Conference on Neural Information Pro-
cessing Systems (NeurIPS)*. 812
813
814
815
816
817

Yanda Chen, Joe Benton, Ansh Radhakrishnan,
Jonathan Uesato, Carson Denison, John Schulman,
Arushi Somani, Peter Hase, Misha Wagner, Fabien
Roger, Vladimir Mikulik, Samuel R. Bowman, Jan
Leike, Jared Kaplan, and Ethan Perez. 2025. [Reason-
ing models don’t always say what they think.](#) *CoRR*,
abs/2505.05410. 818
819
820
821
822
823
824

Micheline T.H. Chi, Miriam Bassok, Matthew W.
Lewis, Peter Reimann, and Robert Glaser. 1989. [Self-
explanations: How students study and use examples
in learning to solve problems.](#) *Cognitive Science*,
13(2):145–182. 825
826
827
828
829

Micheline T.H. Chi, Nicholas De Leeuw, Mei-Hung
Chiu, and Christian Lavancher. 1994. [Eliciting self-
explanations improves understanding.](#) *Cognitive Sci-
ence*, 18(3):439–477. 830
831
832
833

James Chua and Owain Evans. 2025. [Are DeepSeek R1
and other reasoning models more faithful?](#) *CoRR*,
abs/2501.08156. 834
835
836

Yihe Deng, Weitong Zhang, Zixiang Chen, and Quan-
quan Gu. 2023. [Rephrase and respond: Let large
language models ask better questions for themselves.](#)
CoRR, abs/2311.04205. 837
838
839
840

841	Jeremias Lino Ferrao, Ezgi Basar, Khondoker Ittehadul Islam, and Mahrokh Hassani. 2025. What really counts? examining step and token level attribution in multilingual cot reasoning . <i>CoRR</i> , abs/2511.15886.	895
842		896
843		897
844		898
845	Pedro Ferreira, Wilker Aziz, and Ivan Titov. 2025. Truthful or fabricated? using causal attribution to mitigate reward hacking in explanations . <i>CoRR</i> , abs/2504.05294.	899
846		900
847		901
848		902
849	Zhizhang FU, Guangsheng Bao, Hongbo Zhang, Chenkai Hu, and Yue Zhang. 2025. Correlation or causation: Analyzing the causal structures of LLM and LRM reasoning process . <i>CoRR</i> , abs/2509.17380.	903
850		904
851		905
852		906
853	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>CoRR</i> , abs/2305.15717.	907
854		908
855		909
856		910
857	Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. 2024. Training large language models to reason in a continuous latent space . <i>CoRR</i> , abs/2412.06769.	911
858		912
859		913
860		914
861	Peter Hase, Shiyue Zhang, Harry Xie, Mohit Bansal, Percy Liang, Yonatan Bisk, and Dan Roth. 2020. Evaluating explanation faithfulness in natural language inference . In <i>Findings of the Association for Computational Linguistics: EMNLP, 2020</i> . Association for Computational Linguistics.	915
862		916
863		917
864		918
865		919
866		920
867	Minbin Huang, Runhui Huang, Chuanyang Zheng, Jingyao Li, Guoxuan Chen, Han Shi, and Hong Cheng. 2025. Answer-consistent chain-of-thought reinforcement learning for multi-modal large language models . <i>CoRR</i> , abs/2510.10104.	921
868		922
869		923
870		924
871		925
872	Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera Lanham, Daniel M. Ziegler, Tim Maxwell, Newton Cheng, Adam S. Jermyn, Amanda Askill, Ansh Radhakrishnan, Cem Anil, David Duvenaud, Deep Ganguli, Fazl Barez, Jack Clark, Kamal Ndousse, and 20 others. 2024. Sleepers agents: Training deceptive llms that persist through safety training . <i>CoRR</i> , abs/2401.05566.	926
873		927
874		928
875		929
876		930
877		931
878		932
879		933
880		934
881	Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable NLP systems: How should we define and evaluate faithfulness? In <i>Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 4198–4205.	935
882		936
883		937
884		938
885		939
886	Kaiyuan Ji, Yijin Guo, Zicheng Zhang, Xiangyang Zhu, Yuan Tian, Ning Liu, and Guangtao Zhai. 2025. Medomni-45°: A safety-performance benchmark for reasoning-oriented llms in medicine . <i>CoRR</i> , abs/2508.16213.	940
887		941
888		942
889		943
890		944
891	Mengzhao Jia, Zhihan Zhang, Ignacio Cases, Zheyuan Liu, Meng Jiang, and Peng Qi. 2025. Autorubric-1v: Rubric-based generative rewards for faithful multimodal reasoning . <i>CoRR</i> , abs/2510.14738.	945
892		946
893		947
894		948
		949
	Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, and Chelsea Finn. 2025. Counterfactual evaluation of vision-language models for compositional chain-of-thought reasoning . <i>CoRR</i> , abs/2501.04245.	949
		950
	Yubin Kim, Xuhai Xu, Daniel McDuff, and Marzyeh Ghassemi. 2025. Perturbation-based evaluation of visual-language models for medical chain-of-thought reasoning . <i>CoRR</i> , abs/2501.03890.	950
		951
		952
	Aakanksha Kumar, Ranjay Krishna, Aditi Raghunathan, and Percy Liang. 2023. Faithful and efficient reasoning with chain-of-thought distillation . In <i>Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)</i> .	953
		954
	Philippe Laban, Lidiya Murakhovska, Caiming Xiong, and Chien-Sheng Wu. 2024. Are you sure? challenging llms leads to performance drops in the flipflop experiment . <i>CoRR</i> , abs/2311.08596.	955
		956
	Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, and 4 others. 2024. Tulu 3: Pushing frontiers in open language model post-training . <i>CoRR</i> , abs/2411.15124.	957
		958
	Zhenzhong Lan, Junwei Bao, Yujia Qin, Weizhi Wang, and Lei Wang. 2025. Typed chain-of-thought: A framework for verifiable and faithful multi-step reasoning . In <i>Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)</i> .	959
		960
	Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamile Lukosiute, Karina Nguyen, Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Robin Larson, Sam McCandlish, Sandipan Kundu, and 11 others. 2023. Measuring faithfulness in chain-of-thought reasoning . <i>CoRR</i> , abs/2307.13702.	961
		962
	Joshua Ong Jun Leang, Aryo Pradipta Gema, and Shay B. Cohen. 2024. Comat: Chain of mathematically annotated thought improves mathematical reasoning . <i>CoRR</i> , abs/2410.10336.	963
		964
	Jixuan Leng, Cassandra A. Cohen, Zhixian Zhang, Chenyan Xiong, and William W. Cohen. 2025. Semi-structured LLM reasoners can be rigorously audited . <i>CoRR</i> , abs/2505.24217.	965
		966
	Mosh Levy, Zohar Elyoseph, Shauli Ravfogel, and Yoav Goldberg. 2025. State over tokens: Characterizing the role of reasoning tokens . <i>CoRR</i> , abs/2512.12777.	967
		968
	Belinda Z. Li, Zifan Carl Guo, Vincent Huang, Jacob Steinhardt, and Jacob Andreas. 2025a. Training language models to explain their own computations . <i>CoRR</i> , abs/2511.08579.	969
		970

1059	Sanchit Sinha, Oana Frunza, Kashif Rasul, Yuriy Nevmyvaka, and Aidong Zhang. 2025. Chart-rvr: Reinforcement learning with verifiable rewards for explainable chart reasoning . <i>CoRR</i> , abs/2510.10973.	1117
1060		1118
1061		1119
1062		
1063	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, and 431 others. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models . <i>Trans. Mach. Learn. Res.</i> , 2023.	1120
1064		1121
1065		1122
1066		1123
1067		1124
1068		1125
1069		1126
1070		
1071		
1072		
1073	Kaya Stechly, Karthik Valmeekam, Atharva Gundawar, Vardhan Palod, and Subbarao Kambhampati. 2025. Beyond semantics: The unreasonable effectiveness of reasonless intermediate tokens . <i>CoRR</i> , abs/2505.13775.	1127
1074		1128
1075		1129
1076		1130
1077		1131
1078		1132
1079	Qiushi Sun, Zhangyue Yin, Xiang Li, Zhiyong Wu, Xipeng Qiu, and Lingpeng Kong. 2023. Corex: Pushing the boundaries of complex reasoning through multi-model collaboration . <i>CoRR</i> , abs/2310.00280.	1133
1080		1134
1081		1135
1082		
1083	Aadil Syed, Joseph Christopher, Swarnadeep Mishra, Zachary M. Ziegler, Yova Kementchedjhieva, and Stefan Scherer. 2024. Attribution patching outperforms automated circuit discovery . In <i>Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP</i> . Association for Computational Linguistics.	1136
1084		1137
1085		1138
1086		1139
1087		
1088		
1089	Sree Harsha Tanneru, Dan Ley, Chirag Agarwal, and Himabindu Lakkaraju. 2024. On the hardness of faithful chain-of-thought reasoning in large language models . <i>CoRR</i> , abs/2406.10625.	1140
1090		1141
1091		1142
1092		1143
1093	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>Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)</i> .	1144
1094		1145
1095		1146
1096		1147
1097		
1098		
1099	Martin Tutek, Fateme Hashemi Chaleshtori, Ana Marasovic, and Yonatan Belinkov. 2025. Measuring chain of thought faithfulness by unlearning reasoning steps . In <i>Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9946–9971, Suzhou, China. Association for Computational Linguistics.	1148
1100		1149
1101		1150
1102		1151
1103		1152
1104		1153
1105		1154
1106	Prasetya Ajie Utama, Nafise Sadat Moosavi, and Iryna Gurevych. 2025. Cross-lingual generalization of chain-of-thought faithfulness in large language models . <i>CoRR</i> , abs/2501.03245.	1155
1107		1156
1108		1157
1109		1158
1110		1159
1111		1160
1112	Teun van der Weij, Felix Hofstätter, Oliver Jaffe, Samuel F. Brown, and Francis Rhys Ward. 2025. AI sandbagging: Language models can strategically underperform on evaluations . In <i>Proceedings of the International Conference on Learning Representations (ICLR), 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.	1161
1113		1162
1114		1163
1115		1164
1116		1165
		1166
		1167
	Scott Viteri, Max Lamparth, Peter Chatain, and Clark Barrett. 2024. Markovian transformers for informative language modeling . <i>CoRR</i> , abs/2404.18988.	1168
		1169
	Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2023. Interpretability in the wild: a circuit for indirect object identification in GPT-2 small . In <i>Proceedings of the International Conference on Learning Representations (ICLR), 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	1170
		1171
		1172
		1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
		1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282
		1283
		1284
		1285
		1286
		1287
		1288
		1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
		1300
		1301
		1302
		1303
		1304
		1305
		1306
		1307
		1308
		1309
		1310
		1311
		1312
		1313
		1314
		1315
		1316
		1317
		1318
		1319
		1320
		1321
		1322
		1323
		1324
		1325
		1326
		1327
		1328
		1329
		1330
		1331
		1332
		1333
		1334
		1335
		1336
		1337
		1338
		1339
		1340
		1341
		1342
		1343
		1344
		1345
		1346
		1347
		1348
		1349
		1350
		1351
		1352
		1353
		1354
		1355
		1356
		1357
		1358
		1359
		1360
		1361
		1362
		1363
		1364
		1365
		1366
		1367
		1368
		1369
		1370
		1371
		1372
		1373
		1374
		1375
		1376
		1377
		1378
		1379
		1380
		1381
		1382
		1383
		1384
		1385
		1386
		1387
		1388
		1389
		1390
		1391
		1392
		1393
		1394
		1395
		1396
		1397
		1398
		1399
		1400
		1401
		1402
		1403
		1404
		1405
		1406
		1407
		1408
		1409
		1410
		1411
		1412
		1413
		1414
		1415
		1416
		1417
		1418
		1419
		1420
		1421
		1422
		1423
		1424
		1425
		1426
		1427
		1428
		1429
		1430
		1431
		1432
		1433
		1434
		1435
		1436
		1437
		1438
		1439
		1440
		1441
		1442
		1443
		1444
		1445
		1446
		1447
		1448
		1449
		1450
		1451
		1452
		1453
		1454
		1455
		1456
		1457
		1458
		1459
		1460
		1461
		1462
		1463
		1464
		1465
		1466
		1467
		1468
		1469
		1470
		1471
		1472
		1473
		1474
		1475
		1476
		1477
		1478
		1479
		1480
		1481
		1482
		1483
		1484
		1485
		1486
		1487
		1488
		1489
		1490
		1491
		1492
		1493
		1494
		1495
		1496
		1497
		1498
		1499
		1500

1168
1169
1170
1171
1172
1173
1174

1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188

1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205

1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217

A Alignment-Induced Unfaithfulness

Standard alignment techniques (e.g. SFT, RLHF) often induce a form of *stylized unfaithfulness*. Because models are optimized to satisfy human annotators who prefer authoritative and structured explanations, they learn to mimic the *form* of reasoning without maintaining the *substance*.

Superficial Mimicry in SFT. During *Supervised Fine-Tuning* (SFT), studies find that the model primarily captures the tone, reasoning format, and token statistics of the training data, occasionally generating plausible yet meaningless steps (Sinha et al., 2025; Lobo et al., 2025; Hase et al., 2020; Zhang et al., 2025b). This suggests that SFT tends to fit superficial token co-occurrence probabilities rather than learning true causal relationships. Furthermore, cross-domain SFT often compromises accuracy in rigorous fields like logic and mathematics, as the model adopts the linguistic style of the training data without acquiring the necessary underlying logic (Lobo et al., 2025).

Knowledge Distillation. Knowledge distillation is widely employed to transfer reasoning capabilities from large teacher models to smaller, efficient student models. However, its impact on CoT faithfulness is nuanced. While some research suggests that distilling high-quality CoT into small models can improve faithfulness by simplifying the reasoning process (Chua and Evans, 2025), there is a substantial risk of superficial mimicry. Small models often learn to mechanically replicate the teacher’s token sequences without acquiring the underlying causal mechanisms (Gudibande et al., 2023). Consequently, if the teacher’s CoT contains biases or idiosyncratic patterns, distillation amplifies these flaws, trapping the student model in deep pattern matching rather than genuine arithmetic or logical derivation (Lobo et al., 2025).

RLHF and Human-Centric Fabrication. *Reinforcement Learning from Human Feedback* (RLHF) may exacerbate unfaithfulness by aligning models with subjective human preferences rather than truth. Since annotators often favor plausible-sounding explanations over actual faithful ones, models learn to optimize for *persuasion* (Casper et al., 2023). This manifests as *label leakage*: models frequently decide the answer first based on heuristics and then fabricate a rationale to justify it, as humans tend to reward such retrospective consistency (Hase et al., 2020). Similarly, models tend to use an artificially

confident tone to mask logical failures, since humans are more likely to reward confident-sounding reasoning trajectories (FU et al., 2025; Viteri et al., 2024; Ferreira et al., 2025).

RLVR and Reward Hacking. While SFT and RLHF often weaken causal structures (FU et al., 2025), *Reinforcement Learning with Verifiable Rewards* (RLVR) stands out as a more effective method for enforcing ideal logical rigor. Research shows that RLVR effectively resists sycophancy and performs better than SFT in following objective rules (Lambert et al., 2024). However, it introduces new challenges. *Outcome Reward Models* (ORMs), which rely solely on the final answer, can encourage “pseudo-reasoning” that maximizes reward without reflecting the internal process (Huang et al., 2025). In contrast, *Process Reward Models* (PRMs) reward the reasoning process, improving faithfulness in formal tasks like mathematics and logic tasks (Chua and Evans, 2025; Huang et al., 2025). However, PRM still faces a *Verification Gap* in complex, open-ended domains (e.g., Question Answering). Because real-world knowledge is nuanced and ambiguous, it is difficult to build reliable rule-based checkers (like in code or math) to automatically verify the reasoning steps (Min et al., 2023). Thus, while RLVR is a promising paradigm, it currently struggles to guarantee faithfulness outside of formal systems. Moreover, the mechanistic understanding of RLVR is still in the early stages, leaving vast space for future research.

B Pathways for Architectural Decoupling

Recent progress in **Continuous Latent Reasoning** and **Introspective Model Explanation** provides a concrete path toward separating explaining from reasoning.

B.1 Reasoning Beyond Language in Latent Space

Traditional CoT is constrained by the syntax of natural language and the requirements of linear output. In contrast, the *Chain of Continuous Thought* paradigm (Hao et al., 2024) highlights the potential of reasoning within a continuous latent space. COCONUT feeds the last hidden state directly back into the model as a “continuous thought”, bypassing the decoding process into discrete text. This enables the implicit reasoning module to prioritize task performance, mitigating the immediate constraints of syntax or plausibility.

1218
1219
1220
1221

1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248

1249

1250
1251
1252
1253

1254
1255

1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266

1267	In this latent space, models appear to exhibit	<i>oversight-oriented</i> evaluations ask whether inter-	1316
1268	planning capabilities resembling Breadth-First	mediate reasoning exposes early warning signals	1317
1269	Search. Continuous thought vectors may simul-	that enable detection or intervention, prioritizing	1318
1270	taneously encode multiple reasoning paths via su-	monitorability over mechanistic alignment (Chan	1319
1271	perposition, potentially facilitating the pruning of	et al., 2025). <i>Uncertainty-aware</i> metrics quantify	1320
1272	incorrect paths through a process akin to implicit	uncertainty in explanations to distinguish “con-	1321
1273	tree search. This implies that a reasoning module	fidient but unreliable” rationales from calibrated	1322
1274	detached from linguistic constraints could support	ones, but uncertainty alone does not certify faith-	1323
1275	logical structures richer than those strictly bound	fulness (Chen et al., 2024). <i>Human-in-the-loop</i>	1324
1276	by natural language.	judgments support auditing and deployment, yet	1325
1277	B.2 Self-Explaining Modules	they often conflate faithfulness with plausibility	1326
1278	Once the reasoning process becomes implicit, a	and provide limited causal guarantees (Jacovi and	1327
1279	dedicated module is required to translate these la-	Goldberg, 2020). Moreover, <i>external adjudica-</i>	1328
1280	tent states to generate human language. For in-	<i>tion</i> (e.g., graders, checkers, executors) strengthens	1329
1281	stance, Li et al. (2025a) trained an interpreter model	claims about correctness or procedural validity, but	1330
1282	specifically to explain its own internal computa-	remains <i>reliability-oriented</i> unless explicitly linked	1331
1283	tions, observing that such self-interpretation train-	to dependence tests (Kumar et al., 2023). These	1332
1284	ing exhibits remarkable data efficiency.	metrics are most useful when reported alongside	1333
1285	Furthermore, the explanation model demon-	causal evaluations—whether black-box, white-box,	1334
1286	strated the ability to detect prompt bias (see Sec-	or hybrid—as doing so clarifies which specific as-	1335
1287	tion 3.1). This is likely because the interpreter	pect of trustworthiness is actually being measured.	1336
1288	can identify internal signal showing that the model	C.2 Meta-evaluation and benchmark scrutiny	1337
1289	is indicating attention to distracting contexts, and	Meta-evaluation treats evaluation itself as object of	1338
1290	faithfully translate them into linguistic descriptions	study, showing that metrics and benchmarks can	1339
1291	such as “I changed my answer because of context	disagree depending on whether they track plausibil-	1340
1292	X”. This highlights the dual advantages of self-	ity, or behavioral and internal causal dependence	1341
1293	explanation: it not only improves efficiency by	(Radhakrishnan et al., 2025, 2024). This motivates	1342
1294	reusing parts of the model’s computational circuits	reporting multiple complementary metrics (black-	1343
1295	but also achieves greater fidelity through <i>privileged</i>	box, white-box, and hybrid) and explicitly docu-	1344
1296	<i>access</i> to its own internal circuitry.	menting the benchmark’s grounding level, rather	1345
1297	C Extended Discussion on Evaluation	than treating any single benchmark as definitive.	1346
1298	Metrics	D Additional and Peripheral Directions	1347
1299	In appendix, we provide additional details on pe-	In addition to the primary paradigms, studies pur-	1348
1300	ripheral metrics that serve as complementary re-	sues <i>evaluation-driven mitigation</i> , where improved	1349
1301	porting dimensions and discuss meta-evaluation	diagnostics or benchmarks guide model selection	1350
1302	studies that scrutinize existing benchmarks.	and iteration without directly modifying training	1351
1303	C.1 Other peripheral metrics	objectives or inference procedures (Liu et al., 2024).	1352
1304	A small set of peripheral metrics is best viewed as	<i>Human-in-the-loop</i> assessment or correction can	1353
1305	complementary reporting dimensions rather than	improve practical reliability during auditing and	1354
1306	standalone measures of causal faithfulness. In most	deployment, but it scales poorly and provides lim-	1355
1307	cases, they capture adjacent notions of “trustwor-	ited causal guarantees (Jacovi and Goldberg, 2020).	1356
1308	thiness” that relate to CoT faithfulness but do not	Moreover, <i>data-centric</i> heuristics—such as filter-	1357
1309	directly test <i>causal dependence</i> . The key point	ing or curating CoT rationales using surface crite-	1358
1310	is that they primarily measure utility or reliability	ria—can serve as weak regularizers, yet they are	1359
1311	(e.g. monitoring, calibration, human auditability,	best viewed as instances of training- or verification-	1360
1312	external verification), not whether the final deci-	based approaches rather than lone mitigations (Ku-	1361
1313	sion is causally driven by CoT rationales. Thus,	mar et al., 2023).	1362
1314	they should be treated as auxiliary evidence un-		
1315	less paired with causal tests. <i>Monitoring- and</i>		