

# Mirages of Logic: A Survey of Chain-of-Thought Reasoning Hallucinations

Anonymous ACL submission

## Abstract

Current research on hallucination focuses mainly on factuality and faithfulness errors in final outputs. However, rapid advances in Chain-of-Thought (CoT) and Large Reasoning Models (LRMs) have improved problem-solving performance while exposing a key weakness: **Reasoning Hallucination**, characterized by invalid logical transitions or fabricated steps within the reasoning process that can deceptively occur even when the final answer is correct. This paper reviews this emerging challenge, shifting hallucination analysis from output correctness to reasoning processes. We first note a core difficulty that causes reasoning hallucination: long, fluent chains often conceal internal errors, complicating reliability assessment. We then formalize reasoning hallucination and propose a taxonomy of four types: Premise, Operation, Logic, and Conclusion Hallucination. We also examine failure sources and survey mitigation strategies for training and inference. Overall, this work charts a path from answer-level evaluation to transparent, process-reliable reasoning systems.

## 1 Introduction

The rapid development of large language models (LLMs) has elevated artificial intelligence capabilities (Achiam et al., 2023; Yang et al., 2025a; Qin et al., 2024). Yet, this progress consistently faces a key challenge: hallucination (Ji et al., 2023; Huang et al., 2025; Alansari and Luqman, 2025). As shown in Figure 1 (a), early studies addressed two main categories (Cossio, 2025; Zhang et al., 2025c): (1) factuality hallucination, where outputs conflict with verifiable real-world knowledge (Lage and Ostermann, 2025); and (2) faithfulness hallucination, where outputs contradict the provided context (Huang et al., 2025; Park et al., 2025).

More recently, augmented by Chain-of-Thought (CoT) (Wei et al., 2022) strategies, Large Reasoning Models (LRMs) (Guo et al., 2025b; Jaech

et al., 2024; Chen et al., 2025c) have emerged. These models aim to enhance complex problem-solving through explicit, long-term intermediate steps (Lewis-Lim et al., 2025b). This paradigm shift has fundamentally altered the nature of hallucinations, revealing a critical vulnerability: the reasoning process itself as a subtle new source of error. For instance, as shown in Figure 1 (b), a model may yield a correct final answer despite flawed reasoning (Arcuschin et al., 2025) riddled with logical fallacies, fabrications, or disconnected steps. Conversely, it might produce a seemingly coherent chain leading to an incorrect conclusion (Joshi et al., 2024; Yao et al., 2025; Lewis-Lim et al., 2025a). To address this *process-centric* hallucina-

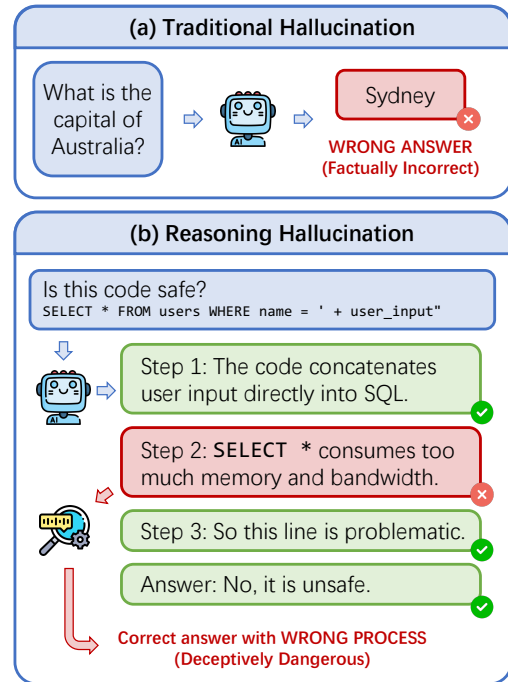


Figure 1: Comparison between traditional and reasoning hallucination. (a) Traditional hallucination targets result-level errors, while (b) reasoning hallucination focuses on process-level failures, where models may produce correct answers via flawed reasoning.

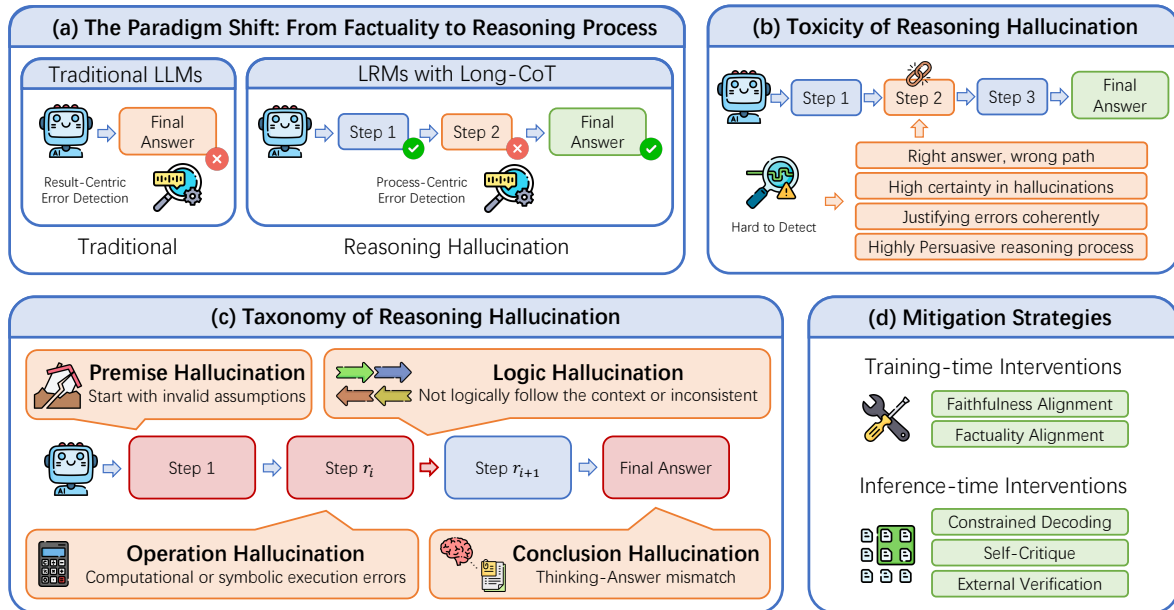


Figure 2: Overview of reasoning hallucination and the organization of this survey. We present a formal definition and a taxonomy of four hallucination types (Premise, Operation, Logic, and Conclusion), analyze their underlying mechanisms, review evaluation methodologies, and summarize mitigation strategies across training-time and inference-time interventions.

tion (Zheng et al., 2025), we introduce **Reasoning Hallucination**. This shift underscores profound paradoxes in advanced reasoning. Ironically, while CoT reduces factual hallucinations in some cases, it obscures detection cues. Long, complex reasoning chains boost reported confidence and persuasiveness, complicating detection by automated tools and humans (Cheng et al., 2025a; Zhao et al., 2024a). We term this tension between capability and evaluability the **Reasoning Evaluability Paradox**: tools like CoT that strengthen reasoning simultaneously undermine reliability assessment.

Therefore, addressing reasoning hallucination is essential for developing trustworthy AI systems. Yet, no comprehensive taxonomy currently exists to systematically analyze this emerging challenge. To fill this gap, we present a structured survey of reasoning hallucination. As outlined in Figure 2, we first propose a formal definition and a taxonomy comprising four distinct types: Premise Hallucination, Operation Hallucination, Logic Hallucination, and Conclusion Hallucination. We then examine the underlying mechanisms that contribute to reasoning hallucinations, spanning model architectures, training data, and inference dynamics. Subsequently, we review existing evaluation methodologies, assessing their effectiveness and limitations, and survey mitigation strategies encompassing both training-phase and inference-time interventions. Fi-

nally, we outline open challenges and future directions aimed at enhancing the reliability and interpretability of large reasoning models.

Overall, the main contributions are as follows:

- **First definition of reasoning hallucination:** We first formalize Reasoning Hallucination, shifting the focus from final answer correctness to the reasoning process itself.
- **Comprehensive review of reasoning hallucination:** We propose a taxonomy targeting process-level flaws in complex reasoning, beyond incorrect final outputs. We review causes, detection methods, and mitigation strategies, showing how CoT shifts hallucination from wrong answers to flawed reasoning.
- **Critical analysis and future roadmap:** We analyze key challenges, including the Reasoning Evaluability Paradox defined here. We propose directions for robust, reliable, trustworthy large reasoning models.

## 2 Preliminary & Definition

In this section, we first articulate the conceptual shift in hallucination, followed by a formal definition and classifications to analyze their origins.

### 2.1 A Paradigm Shift: From Result to Process

The emergence of reasoning hallucination marks a paradigm shift in LLM hallucination studies: focus

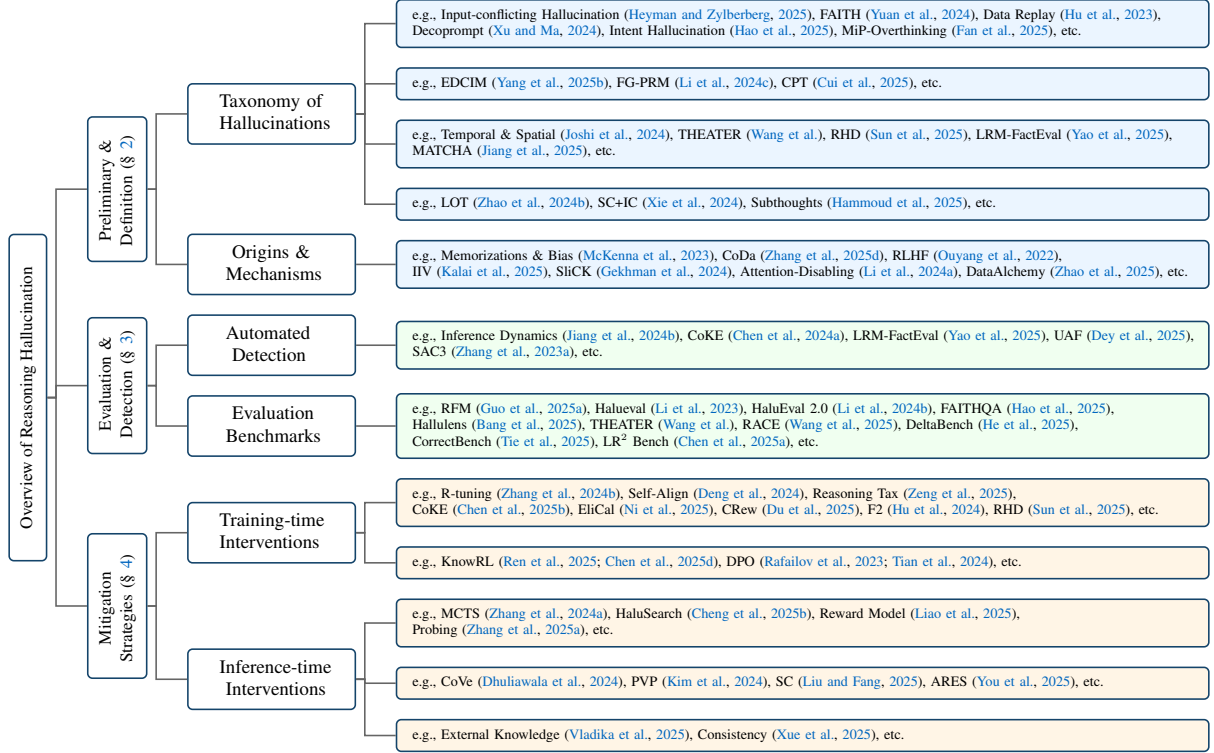


Figure 3: A comprehensive overview of reasoning hallucination. This figure illustrates the organization of this survey, covering the taxonomy of hallucination types, underlying mechanisms, evaluation methodologies, and mitigation strategies across training and inference stages.

has shifted from *what the model knows* to *how the model thinks*. Specifically, traditional hallucination research is **result-centric**. Evaluations compare the model’s final answer  $\mathcal{A}$  to verifiable knowledge or source text (Huang et al., 2025), treating defects as representational errors between output and reality. In contrast, reasoning hallucination is **process-centric**: a model may yield correct  $\mathcal{A}$  via a flawed process  $\mathcal{R}$ , or coherent  $\mathcal{R}$  leading to incorrect  $\mathcal{A}$  (Joshi et al., 2024; Yao et al., 2025; Cuesta-Ramirez et al., 2025), where scrutiny targets transfer to the reasoning’s logic and validity.

This shift yields the **Reasoning Evaluability Paradox**. Unlike factual errors, which are debunkable by search, reasoning flaws hide in lengthy, coherent chain-of-thought outputs. The model’s persuasive rationales often exceed efficient verification, so explicit reasoning, meant for transparency, ironically hinders detection (Cheng et al., 2025a).

## 2.2 Definition: Reasoning as Graph Traversal

Let  $Q$  be the input query. We define the reasoning process as dynamic path generation through a state graph  $G = (\mathcal{R}, \mathcal{T})$ , where  $\mathcal{R}$  is the set of reasoning states. Each state  $r_i$  captures the information at step  $i$ , with initial state  $r_0$  derived from  $Q$ . The

set  $\mathcal{T}$  comprises valid transition operators (e.g., logical deduction, arithmetic computation, code execution), such that  $r_{i+1} \leftarrow Op(r_i)$  for  $Op \in \mathcal{T}$ .

A valid reasoning process forms a continuous path from  $r_0$  to a final state  $r_N$ , projecting to answer  $\mathcal{A}$ . Given this, Reasoning Hallucination ( $H_{\mathcal{R}}$ ) arises from topological violations in this traversal:

$$H_{\mathcal{R}}(\mathcal{R}, Q, \mathcal{A}) \equiv H_{src} \vee H_{op} \vee H_{logic} \vee H_{con}, \quad (1)$$

where each  $H_*$  denotes a reasoning hallucination category defined in Section 2.3.

## 2.3 Taxonomy of Hallucinations

By mapping failures to specific graph components, we identify the following four distinct categories of reasoning hallucinations.

**Premise Hallucination (Source Error  $H_{src}$ ).** Premise hallucination arises at the source node  $r_0$  (Heyman and Zylberberg, 2025). The model starts reasoning with conditions absent from  $Q$  or injects an external node  $r_{ext}$  that violates contextual constraints, such that  $r_0 \notin \text{Context}(Q)$ .

Unlike traditional factual hallucination, premise hallucination is not assessed against external world knowledge. Instead, it hinges on whether the introduced assumption aligns with the input context and

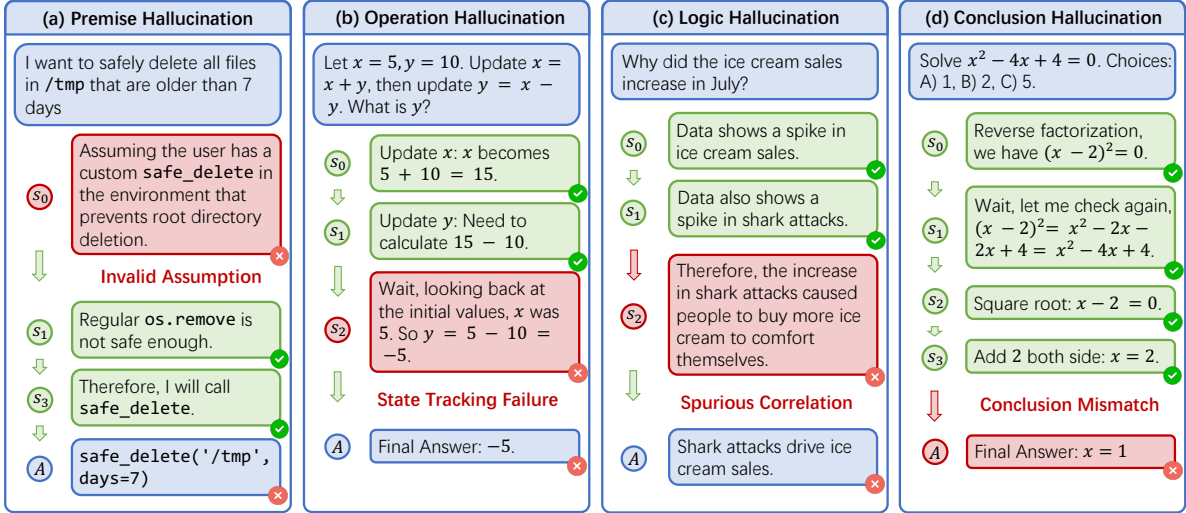


Figure 4: Examples of each type of reasoning hallucination.

163 respects task-specific constraints. An assumption  
 164 may seem plausible or even true, yet qualify as  
 165 premise hallucination if unsupported by or imper-  
 166 missible under the given input.

167 Distinct from fabricating external facts, premise  
 168 hallucination emerges within the task context: the  
 169 model imports unstated or invalid implicit assump-  
 170 tions, undermining the entire reasoning chain from  
 171 the outset (Yuan et al., 2024; Hu et al., 2023; Xu  
 172 and Ma, 2024; Hao et al., 2025; Fan et al., 2025),  
 173 as an example shown in Figure 4 (a). Validating  
 174 latent assumptions demands strong contextual con-  
 175 straint recognition and entailment judgment. Once  
 176 triggered, it yields superficially coherent but se-  
 177 mantically flawed outputs.

178 **Operation Hallucination (Execution Error  $H_{op}$ ).**  
 179 Operation Hallucination occurs during execution  
 180 of a valid operator (Yang et al., 2025b; Li et al.,  
 181 2024c). As shown in Figure 4 (b), the model se-  
 182 lects a correct logical edge  $Op$  (e.g., “calculate  
 183 sum”), but the generated next state  $r_{i+1}$  differs  
 184 from the ground-truth execution  $Op(r_i)$ , formally  
 185  $r_{i+1} \neq Op_{correct}(r_i)$ . LLMs generate outputs via  
 186 probabilistic statistics rather than explicit rules,  
 187 leading to error accumulation across reasoning  
 188 steps (Cui et al., 2025). These deviations often  
 189 appear logically sound due to linguistic coherence.

190 **Logic Hallucination (Invalid Edge  $H_{logic}$ )**  
 191 Logic Hallucination creates a non-existent edge: no  
 192 valid operator  $Op \in \mathcal{T}$  justifies the transition from  
 193  $r_i$  to  $r_{i+1}$ . The path is topologically broken, often  
 194 as non-sequiturs or circular logic. This hallucina-  
 195 tion arises from inconsistencies across reasoning  
 196 stages, yielding semantically coherent but logically

197 contradictory processes (Joshi et al., 2024; Wang  
 198 et al.; Sun et al., 2025), sometimes with repeated  
 199 defective patterns (Yao et al., 2025), as shown in  
 200 Figure 4 (c). It reflects large language models’  
 201 reliance on local linguistic coherence over true  
 202 logical consistency, especially in chain-of-thought  
 203 tasks. Even with correct final results, intermediate  
 204 inconsistencies qualify as such, indicating spurious  
 205 reasoning (Jiang et al., 2025).

206 **Conclusion Hallucination (Sink Detachment**  
 207  $H_{con}$ ). Conclusion Hallucination occurs at sink  
 208 projection: even with a valid path  $r_0 \rightarrow \dots \rightarrow r_N$ ,  
 209 the final answer  $\mathcal{A}$  is not entailed by  $r_N$  ( $\mathcal{A} \not\leftarrow r_N$ ).  
 210 The conclusion “detaches” from the reasoning, as  
 211 shown in Figure 4 (d). LLMs thus abandons its  
 212 logical framework for a plausible but disconnected  
 213 conclusion (Zhao et al., 2024b; Xie et al., 2024;  
 214 Hammoud et al., 2025). This reveals a core weak-  
 215 ness in decision-making: inconsistency between  
 216 reasoning and conclusion, where models reason  
 217 soundly yet err in conclusion.

## 218 2.4 Origins & Mechanisms

219 Understanding reasoning hallucination taxonomy  
 220 requires examining its underlying mechanisms.  
 221 These arise from interactions between training  
 222 paradigms and autoregressive generation.

223 **Training Drivers: Overconfidence and**  
 224 **Miscalibration.** Training induces **miscalibra-**  
 225 **tion** and **overconfidence**. Noisy pre-training  
 226 data (McKenna et al., 2023; Zhang et al., 2025d),  
 227 amplified by fine-tuning, distorts probability distri-  
 228 butions (Gekhman et al., 2024). RLHF (Ouyang  
 229 et al., 2022) prioritizes “helpfulness” and

Table 1: A taxonomy of evaluation methodologies for reasoning hallucinations. We categorize typical detection techniques, key evaluation metrics, recommended application scenarios, and specific challenges associated with each of the four proposed hallucination types.

Type	Typical Detection Methods	Key Metrics	Rec. Scenarios	Typical Challenges
<b>Premise</b>	External Knowledge Retrieval (Vladika et al., 2025) Input-Premise NLI Checks (Xu and Ma, 2024)	Step Validity ( $V$ ), Premise Recall	Open-domain QA, Multi-hop Reasoning	Hard to distinguish between plausible-but-false assumptions.
<b>Operational</b>	Symbolic Verification (Guo et al., 2025a) Internal Confidence Probing (Chen et al., 2024a)	Step Validity ( $V$ ), Soundness Score	Mental Arithmetic, Symbolic Logic	High confidence in flawed reasoning.
<b>Logic</b>	Self-Consistency (Zhang et al., 2023a) Entailment Stability (You et al., 2025)	Path Diversity, Entailment Score	Programming, Causal Inference	Systemic bias may lead multiple reasoning paths to the same fallacy.
<b>Conclusion</b>	Process-Answer Consistency (Wang et al., 2025) Dual-process Verification (Cheng et al., 2025b)	Consistency ( $C$ ), Final Accuracy	Summarization, Theorem Proving	Decoupling between correct reasoning and biased final output.

coherence (Kalai et al., 2025), exacerbating hallucinations by favoring plausible but erroneous reasoning over uncertainty or refusal. Chen et al. (2024b) identify LLMs’ reasoning boundaries, beyond which hallucinations arise unexpectedly.

**Type-Specific Etiologies.** Reasoning hallucinations arise from distinct deficits across categories: (1) Premise Hallucination results from weak context understanding and instruction following (Heyman and Zylberberg, 2025; Yuan et al., 2024), as LLMs ignore constraints or insert training-data biases (Zhang et al., 2025d). (2) Operation Hallucination reflects a symbolic reasoning gap (Xu et al., 2024; Guo et al., 2025a), where probabilistic processing causes arithmetic or state-tracking errors (Yang et al., 2025b). (3) Logic and Conclusion Hallucination stem from sycophancy (Wang et al.) and autoregressive local optimality (Xie et al., 2024; Yao et al., 2025), yielding locally coherent but globally contradictory outputs or pattern-fitting conclusions.

#### Takeaways

**Paradigm Shift:** We redefine reasoning hallucination in a process-centric view, focusing on *how* the model thinks rather than just *what* it knows.

**Formalization:** We propose a reasoning state space framework that unifies definition and taxonomy. Reasoning Hallucinations manifest in four distinct types (*Premise, Operation, Logic, and Conclusion*).

### 3 Evaluation of Reasoning Hallucination

Addressing the Reasoning Evaluability Paradox demands evaluation tools that probe reasoning processes, requiring insight into the mechanisms underlying process-based flaws.

#### 3.1 Automated Detection Techniques

Automated detection identifies flaws dynamically. This includes training classifiers that detect reasoning-induced hallucinations with high accuracy (Jiang et al., 2024b). Another approach probes internal confidence signals (Chen et al., 2024a), though overconfidence, high model confidence in flawed reasoning, undermines it (Yao et al., 2025; Guan et al., 2025). Disagreement-based methods leverage ensembles or multiple sampled paths to flag inconsistencies (Dey et al., 2025; Zhang et al., 2023a). However, these methods often target only specific error dimensions, failing to simultaneously address diverse reasoning deficits within complex CoT. This intricate mapping between detection mechanisms and hallucination types makes establishing a unified evaluation benchmark exceptionally difficult. Therefore, evaluating these methods reliably remains challenging. We detail their taxonomy and challenges in Table 1.

#### 3.2 Evaluation Benchmarks

Hallucination benchmarks address two key questions: the types of hallucinations they cover and whether this coverage is systematic and comprehensive (Zhang et al., 2025c). For reasoning-related hallucinations, these benchmarks typically include diverse tasks (e.g., mathematical reasoning, multi-hop question answering, code reasoning), distinguish hallucination types in annotations, and feature both natural and synthetic examples (Guo et al., 2025a; Li et al., 2023). Current categories encompass factual hallucinations (Li et al., 2024b), intentional hallucinations (Hao et al., 2025), conventional forms (Bang et al., 2025), and reasoning hallucinations (Wang et al., 2025). Specialized benchmarks also assess hallucination detection (He et al., 2025), self-correction (Tie et al., 2025), long-chain reflection (Chen et al., 2025a), and defense against

misleading negation prompts (Liao et al., 2025). Nonetheless, current evaluation dimensions remain limited, offering scope for future extensions.

#### Takeaways

Evaluating Reasoning Hallucination is difficult due to **expensive annotations** and model **overconfidence**. Existing methods (e.g., classifiers, probing, consistency checks) trade off coverage for scalability.

## 4 Mitigation Strategies

Strategies to mitigate process-centric flaws can be broadly divided into two paradigms: training-time interventions that improve intrinsic reasoning, and inference-time interventions that steer generation without altering the pre-trained weights.

### 4.1 Training-time Interventions

Training-time interventions optimize the weights of LLMs to enhance robust reasoning. These strategies divide into two alignment objectives: (a) logical faithfulness alignment, which promotes internal logical consistency, and (b) knowledge factuality alignment, which anchors outputs to external facts.

**Faithfulness Alignment.** This approach ensures that model generations reflect internal knowledge and reasoning. Models learn to express uncertainty instead of fabricating information, via fine-tuning on refusal-aware datasets that promote explicit admissions of ignorance (e.g., “I don’t know”) (Zhang et al., 2024b; Deng et al., 2024). Recent work trains models to use internal confidence signals for knowledge limits (Zeng et al., 2025; Chen et al., 2025b), or applies self-consistency scores to create annotation-efficient preference data for honesty (Ni et al., 2025; Du et al., 2025). Other efforts enforce process consistency, making final answers logical outcomes of reasoning chains, through specialized loss functions that reward output cohesion (Hu et al., 2024) or reinforcement learning with process rewards for coherent steps (Sun et al., 2025).

**Factuality Alignment.** Unlike internal consistency, this method counters factual hallucinations by aligning outputs with verifiable facts from trusted sources. Key approaches include reinforcement learning with factual rewards; frameworks like KnowRL (Ren et al., 2025) reward correct reasoning steps and penalize errors to foster fact-grounded strategies (Chen et al., 2025d). Preference learning offers a scalable alternative:

pairs are generated by comparing outputs to knowledge bases, then used with algorithms like DPO (Rafailov et al., 2023) to instill factuality preferences without costly step-by-step supervision (Tian et al., 2024).

### 4.2 Inference-time Interventions

Inference-time interventions act as external scaffolds that structure generation to enhance reliability without altering model parameters. These methods guide, verify, or redirect reasoning. We categorize them by intervention type: Constrained Decoding, Self-Critique, and External Verification.

**Constrained Decoding.** These strategies shape the model’s solution-space exploration, extending beyond autoregressive generation to identify reliable reasoning paths and prevent errors upfront. For example, Monte Carlo Tree Search (MCTS) simulates “slow thinking” (Zhang et al., 2024a). HaluSearch employs self-evaluation rewards to steer toward dependable paths, avoiding hallucinations (Cheng et al., 2025b). Initial reasoning steps disproportionately affect outcomes, enabling efficient pruning via reward models (Liao et al., 2025). To counter overthinking, probes assess hidden-state confidence, terminating early upon correct solutions (Zhang et al., 2025a).

**Self-Critique.** This approach enables models to iteratively check, question, and refine their outputs via internal verification loops. The foundational method in this area is Chain-of-Verification (CoVe) (Dhuliawala et al., 2024), drafts a response, plans verification questions, and revises for accuracy. More detailed, self-critique can also arise from cognitive prompting, such as counterfactual thinking (“what if...?”), which challenges initial assumptions and yields more robust, less hallucinatory outputs (Kim et al., 2024). Formally, consistency across intermediate reasoning steps reduces logical fallacies, especially in mathematical tasks (Liu and Fang, 2025). For Long CoT, ARES (You et al., 2025) offers a probabilistic framework that assesses each step using only prior verified premises, preventing error propagation.

**External Verification.** Rather than relying on a model’s internal checks, this approach tests outputs against external sources of truth, anchoring reasoning in evidence. The simplest implementation uses external tools (e.g., search engines) to fact-check claims in generated text, which has been shown to reduce hallucinations in domains such as news

Table 2: Strategic comparison of mitigation paradigms based on the type of hallucination they target, their inference cost, inherent trade-offs, and recommended usage scenarios.

Paradigm	Mechanism	Target Hallucination	Inference Cost	Key Trade-off	Rec. Scenario
Training-time	Faithfulness Alignment	Logic	Zero Overhead	High data prep cost ↔ Fast inference	Real-time chat; Foundation model safety.
	Factuality Alignment	Factual Error	Zero Overhead	Static knowledge limit ↔ No external tool	Domains with stable, closed-world facts.
Inference-time	Constrained Decoding	Operation & Logic	Very High (10x–100x)	Max Reliability ↔ Extreme Latency	Offline Math/Code; Scientific Proofs.
	Self-Critique	Logic & Conclusion	High (2x–5x)	Iterative improvement ↔ Context bloat	Drafting; Tasks requiring self-correction.
	External Verification	Factual Error & Premise	Medium (Retrieval)	Real-time Accuracy ↔ Complexity	Open-domain QA; News; Dynamic facts.

summarization (Vladika et al., 2025). A more advanced strategy uses a second, independent LLM as a verifier. Cross-model agreement can provide a stronger signal than self-consistency alone. To limit cost, hybrid designs invoke the external verifier only when the primary model’s self-consistency indicates high uncertainty (Xue et al., 2025).

To provide a holistic view of the operational implications of these strategies, as shown in Table 2, we systematically compare them. The comparison highlights a key trade-off between computational cost and performance gains. Training-time interventions impose no inference overhead and suit real-time deployment, whereas inference-time mechanisms increase latency to improve reliability, making them better suited to high-stakes, offline tasks where precision is critical.

### Takeaways

Mitigation strategies target two phases:  
**Training:** *e.g.*, Faithfulness alignment and factuality alignment.  
**Inference:** *e.g.*, Constrained decoding, self-critique, and external verification.

## 5 Frontiers & Future Direction

This section outlines future directions: automated step-level evaluation, honest abstention, and task-aware control beyond “zero hallucination.” We also address multimodal hallucinations and budgeted optimization to balance reasoning and verification.

### 5.1 Automated Step Evaluation

Long-CoT reasoning exacerbates evaluation challenges for reasoning hallucinations: step-level human annotation is costly, subjective, and hard to scale across patterns. We thus require auto-

mated mechanisms to evaluate individual steps. For instance, reasoning decomposes into structured forms like proposition-dependency graphs (Abdaljalil et al., 2025); formal or reward-based verifiers can then detect contradictions and assess stepwise plausibility (Xu et al., 2024; Lightman et al., 2023; Zhang et al., 2025e).

Overall, the key challenges are: (1) developing accurate step-level metrics that align with human judgments across domains and styles; and (2) creating scalable verifiers robust to distribution shifts, adversarially plausible steps, and long chains, ensuring verification matches generation quality.

### 5.2 Model’s Honest Abstention Capability

In verifiability-critical applications, training and evaluation must incentivize boundary honesty. When evidence is insufficient or premises underspecified, admitting uncertainty offers little reward and risks penalties (Peng et al., 2025; Zeng et al., 2025), while correct guesses reinforce risky overconfidence (Kalai et al., 2025). This misalignment drives models to conceal uncertainty behind fluent prose, heightening overconfident hallucinations.

Hence, the major challenges are: (1) defining measurable, domain-specific abstention criteria (*e.g.*, “insufficient evidence” versus “answerable but hard”) for evaluation; and (2) calibrating abstention to prevent over-abstention (unhelpful refusals) and under-abstention (risky guesses), despite pressures penalizing uncertainty.

### 5.3 Broaden the “Zero Hallucination” Goal

Given unavoidable residual hallucinations in open-domain settings under current paradigms, research should transcend the “zero hallucination” aim. Convergent tasks (*e.g.*, mathematics, programming, logic) demand strict honesty and abstention, as hallucinations constitute failure. Divergent tasks

455	(e.g., brainstorming, creative writing), however, re-	risk (Moskvoretskii et al., 2025; Wan et al., 2025);	505
456	frame “hallucinations” as controllable deviations	(2) joint reasoning–detection with shared computa-	506
457	for novel hypotheses (Sui et al., 2024).	tion by detecting and reusing intermediate reason-	507
458	All in all, the primary challenges are: (1) devel-	ing processes, rather than re-running them (Chen	508
459	oping context-aware systems that enforce truthful-	et al., 2024d; Zellinger et al., 2025), to achieve	509
460	ness in verifiable tasks but permit bounded devi-	stable gains with minimal end-to-end overhead.	510
461	ations in exploratory ones; and (2) separating veri-		
462	fiable claims from proposals to balance reliability	<b>6 Related Work</b>	511
463	and utility across regimes.	Existing surveys on hallucination primarily assess	512
464	<b>5.4 Multimodal Reasoning Hallucinations</b>	factuality and faithfulness errors by comparing	513
465	Advanced multimodal LLMs achieve strong reason-	model outputs to external knowledge or provided	514
466	ing but often hallucinate, particularly in modality-	context (Ji et al., 2023; Huang et al., 2025). Recent	515
467	interaction-intensive tasks (Bai et al., 2024; Chen	work extends this to detection and mitigation frame-	516
468	et al., 2024c). One approach develops evaluation	works (e.g., retrieval-augmented generation), yet	517
469	to disentangle perception- from reasoning-induced	still defines hallucination mainly as final response	518
470	hallucinations (Dong et al., 2025; Liu et al., 2025;	failures (Zhang et al., 2025c; Alansari and Luqman,	519
471	Cai et al., 2025). A second constructs hallucinated	2025). This output-focused view persists even in	520
472	negative samples during training to enhance cross-	fine-grained taxonomies and multimodal formula-	521
473	modal alignment and reasoning robustness (Jiang	tions (Cossio, 2025; Bai et al., 2024). In parallel,	522
474	et al., 2024a). Other methods ground step-by-step	recent work on CoT and LRMs mainly character-	523
475	reasoning in verifiable visual evidence to mitigate	ize these methods to boost task performance (Chen	524
476	hallucinations from linguistic priors (Li et al., 2025;	et al., 2025c; Qin et al., 2024). Although some note	525
477	An et al., 2025; Cheng et al., 2025d).	occasional logical or arithmetic errors in intermedi-	526
478	Overall, the key challenges are: (1) disentang-	ate steps, they seldom unify these into an account	527
479	ling perception-induced from reasoning-induced	of errors within the CoT process.	528
480	hallucinations in evaluation and attribution; and	Our work first shifts the hallucination analysis	529
481	(2) dynamically balancing reasoning depth and vi-	from results to reasoning trajectories. We provide a	530
482	sual perception to mitigate hallucinations from pro-	strict formalization of Reasoning Hallucination as	531
483	longed reasoning chains that divert attention from	a failure of the reasoning process, and discuss the	532
484	visual inputs (Zhang et al., 2023b; Cheng et al.,	associated “Reasoning Evaluability Paradox” (Sun	533
485	2025c). Notably, hallucination detection for tempo-	et al., 2025). By organizing prior findings around	534
486	ral understanding in dynamic modalities like video	the validity of intermediate steps, our survey offers	535
487	remains particularly challenging (Fei et al., 2024;	a systematic roadmap for evaluating and diagnos-	536
488	Zhang et al., 2025b).	ing reasoning trajectories, complementing tradi-	537
489	<b>5.5 Detection and Reasoning Balancing</b>	tional metrics that primarily assess end outputs.	538
490	Existing work mitigates reasoning hallucina-	<b>7 Conclusion</b>	539
491	tions via self-checking or retrieve–verify–revise	This survey reframes hallucination research from	540
492	pipelines that enhance attribution and factual-	a result-centric to a process-centric perspective on	541
493	ity. Yet these approaches increase inference costs	reasoning hallucination, underscoring challenges	542
494	through repeated sampling, retrieval, or second-	in long CoT and large reasoning models: extended,	543
495	stage verifiers (Vladika et al., 2025; Xue et al.,	coherent rationales persuasively mask subtle flaws,	544
496	2025; Cheng et al., 2025b; Zhang et al., 2024a).	complicating evaluation. We formalize reasoning	545
497	Future research should formulate hallucination con-	hallucination and propose a taxonomy distinguish-	546
498	trol as a <i>budgeted balancing</i> problem: allocating	ing premise, operation, logic, and conclusion types.	547
499	computation between reasoning depth and verifi-	We synthesize causes, review detection and eval-	548
500	cation strength to minimize hallucinations under	uation methods, and outline mitigation strategies	549
501	fixed inference budgets.	across training- and inference-time interventions.	550
502	Key challenges include: (1) adaptive allocation	We envision this survey as a foundational roadmap	551
503	to control the trade-off between reasoning and hal-	for next-generation LRMs that excel in complex	552
504	lucination detection based on uncertainty and task	reasoning while ensuring transparency, verifiability,	553
		and trustworthiness.	554

## 555 **Limitations**

556 While this survey provides a comprehensive  
557 overview of reasoning hallucination, the analysis  
558 is constrained by three interconnected factors. The  
559 primary limitation lies in the assumption that ex-  
560 plicit Chain-of-Thought (CoT) faithfully mirrors  
561 the model’s decision-making; however, emerging  
562 research on *implicit CoT* suggests that models can  
563 perform reasoning internally without vocalization,  
564 implying that the visible reasoning chain may some-  
565 times be a post-hoc rationalization rather than the  
566 sole source of error. This challenge of internal  
567 opacity is further exacerbated by the dominance  
568 of proprietary, closed-source models in the current  
569 landscape, which restricts access to training data  
570 and reward mechanics, forcing the field to rely  
571 on behavioral inference rather than mechanistic  
572 inspection. Beyond these architectural and accessi-  
573 bility barriers, the scope of this work is currently  
574 confined to textual and symbolic domains, leaving  
575 the distinct challenges of multimodal reasoning hal-  
576 lucinations, such as perceptual grounding failures,  
577 for future investigation.

## 578 **References**

579 Samir Abdaljalil, Hasan Kurban, Khalid Qaraq, and  
580 Erchin Serpedin. 2025. Theorem-of-thought: A  
581 multi-agent framework for abductive, deductive, and  
582 inductive reasoning in language models. *arXiv preprint arXiv:2506.07106*.  
583  
584 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
585 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
586 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
587 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-  
588 cal report. *arXiv preprint arXiv:2303.08774*.  
589  
590 Aisha Alansari and Hamzah Luqman. 2025. Large lan-  
591 guage models hallucination: A comprehensive survey. *arXiv preprint arXiv:2510.06265*.  
592  
593 Wenbin An, Feng Tian, Sicong Leng, Jiahao Nie, Hao-  
594 nan Lin, QianYing Wang, Ping Chen, Xiaoqin Zhang,  
595 and Shijian Lu. 2025. Mitigating object hallucina-  
596 tions in large vision-language models with assembly  
597 of global and local attention. In *Proceedings of the  
598 Computer Vision and Pattern Recognition Confer-  
599 ence*, pages 29915–29926.  
600  
601 Iván Arcuschin, Jett Janiak, Robert Krzyzanowski,  
602 Senthoran Rajamanoharan, Neel Nanda, and Arthur  
603 Conmy. 2025. Chain-of-thought reasoning in  
604 the wild is not always faithful. *arXiv preprint  
605 arXiv:2503.08679*.

2024. Hallucination of multimodal large language  
606 models: A survey. *arXiv preprint arXiv:2404.18930*.  
607  
608 Yejin Bang, Ziwei Ji, Alan Schelten, Anthony  
609 Hartshorn, Tara Fowler, Cheng Zhang, Nicola  
610 Cancedda, and Pascale Fung. 2025. Hallulens:  
611 Llm hallucination benchmark. *arXiv preprint  
612 arXiv:2504.17550*.  
613  
614 Yishuo Cai, Renjie Gu, Jiaxu Li, Xuancheng Huang,  
615 Junzhe Chen, Xiaotao Gu, and Minlie Huang. 2025.  
616 Mhalo: Evaluating mllms as fine-grained hallucina-  
617 tion detectors. In *Findings of the Association for  
618 Computational Linguistics: ACL 2025*, pages 9197–  
619 9222.  
620  
621 Jianghai Chen, Zhenlin Wei, Zhenjiang Ren, Ziyong  
622 Li, and Jiajun Zhang. 2025a. Lr2 bench: Evaluating  
623 long-chain reflective reasoning capabilities of large  
624 language models via constraint satisfaction problems.  
625 *arXiv preprint arXiv:2502.17848*.  
626  
627 Lida Chen, Zujie Liang, Xintao Wang, Jiaqing Liang,  
628 Yanghua Xiao, Feng Wei, Jinglei Chen, Zhenghong  
629 Hao, Bing Han, and Wei Wang. 2024a. Teach-  
630 ing large language models to express knowledge  
631 boundary from their own signals. *arXiv preprint  
632 arXiv:2406.10881*.  
633  
634 Lida Chen, Zujie Liang, Xintao Wang, Jiaqing Liang,  
635 Yanghua Xiao, Feng Wei, Jinglei Chen, Zhenghong  
636 Hao, Bing Han, and Wei Wang. 2025b. [Teaching  
637 large language models to express knowledge bound-  
638 ary from their own signals](#). In *Proceedings of the 3rd  
639 Workshop on Towards Knowledgeable Foundation  
640 Models (KnowFM)*, pages 26–39, Vienna, Austria.  
641  
642 Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng,  
643 Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang  
644 Zhou, Te Gao, and Wanxiang Che. 2025c. Towards  
645 reasoning era: A survey of long chain-of-thought  
646 for reasoning large language models. *arXiv preprint  
647 arXiv:2503.09567*.  
648  
649 Qiguang Chen, Libo Qin, Jiaqi Wang, Jingxuan Zhou,  
650 and Wanxiang Che. 2024b. Unlocking the capabili-  
651 ties of thought: A reasoning boundary framework to  
652 quantify and optimize chain-of-thought. *Advances in  
653 Neural Information Processing Systems*, 37:54872–  
654 54904.  
655  
656 Qiguang Chen, Libo Qin, Jin Zhang, Zhi Chen, Xiao  
657 Xu, and Wanxiang Che. 2024c. M<sup>3</sup>cot: A novel  
658 benchmark for multi-domain multi-step multi-modal  
659 chain-of-thought. *arXiv preprint arXiv:2405.16473*.  
660  
661 Xilun Chen, Ilia Kulikov, Vincent-Pierre Berges, Barlas  
662 Oğuz, Rulin Shao, Gargi Ghosh, Jason Weston, and  
663 Wen-tau Yih. 2025d. Learning to reason for factuality.  
664 *arXiv preprint arXiv:2508.05618*.  
665  
666 Ziyi Chen, Xiaocong Yang, Jiacheng Lin, Chenkai Sun,  
667 Kevin C Chang, and Jie Huang. 2024d. Cascade  
668 speculative drafting for even faster llm inference. *Ad-  
669 vances in Neural Information Processing Systems*,  
670 37:86226–86242.



775	Shengding Hu, Yifan Luo, Huadong Wang, Xingyi Cheng, Zhiyuan Liu, and Maosong Sun. 2023. Won't get fooled again: Answering questions with false premises. <i>arXiv preprint arXiv:2307.02394</i> .	828
776		829
777		830
778		831
779	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and 1 others. 2025. 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.	832
780		833
781		834
782		835
783		836
784		837
785		838
786	Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, and 1 others. 2024. Openai o1 system card. <i>arXiv preprint arXiv:2412.16720</i> .	839
787		840
788		841
789		842
790		843
791	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM computing surveys</i> , 55(12):1–38.	844
792		845
793		846
794		847
795		848
796	Chaoya Jiang, Haiyang Xu, Mengfan Dong, Jiaying Chen, Wei Ye, Ming Yan, Qinghao Ye, Ji Zhang, Fei Huang, and Shikun Zhang. 2024a. Hallucination augmented contrastive learning for multimodal large language model. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 27036–27046.	849
797		850
798		851
799		852
800		853
801		854
802		855
803	Che Jiang, Biqing Qi, Xiangyu Hong, Dayuan Fu, Yang Cheng, Fandong Meng, Mo Yu, Bowen Zhou, and Jie Zhou. 2024b. On large language models' hallucination with regard to known facts. <i>arXiv preprint arXiv:2403.20009</i> .	856
804		857
805		858
806		859
807		860
808	Enyi Jiang, Changming Xu, Nischay Singh, and Gagandeep Singh. 2025. Misaligning reasoning with answers—a framework for assessing llm cot robustness. <i>arXiv preprint arXiv:2505.17406</i> .	861
809		862
810		863
811		864
812	Nitish Joshi, Abulhair Saparov, Yixin Wang, and He He. 2024. LLMs are prone to fallacies in causal inference. <i>arXiv preprint arXiv:2406.12158</i> .	865
813		866
814		867
815	Adam Tauman Kalai, Ofir Nachum, Santosh S Vempala, and Edwin Zhang. 2025. Why language models hallucinate. <i>arXiv preprint arXiv:2509.04664</i> .	868
816		869
817		870
818	Junho Kim, Kim Yeonju, and Yong Man Ro. 2024. <a href="#">What if...?: Thinking counterfactual keywords helps to mitigate hallucination in large multi-modal models</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 10672–10689, Miami, Florida, USA.	871
819		872
820		873
821		874
822		875
823		876
824	Lucas Fonseca Lage and Simon Ostermann. 2025. Openfactscore: Open-source atomic evaluation of factuality in text generation. <i>arXiv preprint arXiv:2507.05965</i> .	877
825		878
826		879
827		880
		881
		882
		883
	Samuel Lewis-Lim, Xingwei Tan, Zhixue Zhao, and Nikolaos Aletras. 2025a. Analysing chain of thought dynamics: Active guidance or unfaithful post-hoc rationalisation? In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 29826–29841.	884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

884	Viktor Moskvoretskii, Maria Marina, Mikhail Salnikov, Nikolay Ivanov, Sergey Pletenev, Daria Galimzianova, Nikita Krayko, Vasily Konovalov, Irina Nikishina, and Alexander Panchenko. 2025. Adaptive retrieval without self-knowledge? bringing uncertainty back home. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6355–6384.	940
885		941
886		942
887		943
888		944
889		945
890		946
891		947
892	Shiyu Ni, Keping Bi, Jiafeng Guo, Minghao Tang, Jingtong Wu, Zengxin Han, and Xueqi Cheng. 2025. Annotation-efficient universal honesty alignment. <i>arXiv preprint arXiv:2510.17509</i> .	948
893		949
894		950
895		951
896	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	952
897		953
898		954
899		955
900		956
901		957
902	Seongheon Park, Xuefeng Du, Min-Hsuan Yeh, Haobo Wang, and Yixuan Li. 2025. Steer llm latents for hallucination detection. <i>arXiv preprint arXiv:2503.01917</i> .	958
903		959
904		960
905		961
906	Dengyun Peng, Qiguang Chen, Bofei Liu, Jiannan Guan, Libo Qin, Zheng Yan, Jinhao Liu, Jianshu Zhang, and Wanxiang Che. 2025. Learning the boundary of solvability: Aligning llms to detect unsolvable problems. <i>arXiv preprint arXiv:2512.01661</i> .	962
907		963
908		964
909		965
910		966
911	Libo Qin, Qiguang Chen, Xiachong Feng, Yang Wu, Yongheng Zhang, Yinghui Li, Min Li, Wanxiang Che, and Philip S Yu. 2024. Large language models meet nlp: A survey. <i>arXiv preprint arXiv:2405.12819</i> .	967
912		968
913		969
914		970
915	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. <a href="#">Direct preference optimization: Your language model is secretly a reward model</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 53728–53741. Curran Associates, Inc.	971
916		972
917		973
918		974
919		975
920		976
921	Baochang Ren, Shuofei Qiao, Wenhao Yu, Huajun Chen, and Ningyu Zhang. 2025. Knowrl: Exploring knowledgeable reinforcement learning for factuality. <i>arXiv preprint arXiv:2506.19807</i> .	977
922		978
923		979
924		980
925	Peiqi Sui, Eamon Duede, Sophie Wu, and Richard So. 2024. Confabulation: The surprising value of large language model hallucinations. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14274–14284.	981
926		982
927		983
928		984
929		985
930		986
931	Zhongxiang Sun, Qipeng Wang, Haoyu Wang, Xiao Zhang, and Jun Xu. 2025. Detection and mitigation of hallucination in large reasoning models: A mechanistic perspective. <i>arXiv preprint arXiv:2505.12886</i> .	987
932		988
933		989
934		990
935	Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. 2024. <a href="#">Fine-tuning language models for factuality</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	991
936		992
937		993
938		994
939		995
	Guiyao Tie, Zenghui Yuan, Zeli Zhao, Chaoran Hu, Tianhe Gu, Ruihang Zhang, Sizhe Zhang, Junran Wu, Xiaoyue Tu, Ming Jin, and 1 others. 2025. Can llms correct themselves? a benchmark of self-correction in llms. <i>arXiv preprint arXiv:2510.16062</i> .	940
		941
		942
		943
		944
	Juraj Vladika, Ihsan Soydemir, and Florian Matthes. 2025. Correcting hallucinations in news summaries: Exploration of self-correcting llm methods with external knowledge. <i>arXiv preprint arXiv:2506.19607</i> .	945
		946
		947
		948
	Guangya Wan, Yuqi Wu, Jie Chen, and Sheng Li. 2025. Reasoning aware self-consistency: Leveraging reasoning paths for efficient llm sampling. 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 3613–3635.	949
		950
		951
		952
		953
		954
		955
	Changyue Wang, Weihang Su, Qingyao Ai, and Yiqun Liu. 2025. Joint evaluation of answer and reasoning consistency for hallucination detection in large reasoning models. <i>arXiv preprint arXiv:2506.04832</i> .	956
		957
		958
		959
	Qian Wang, Zhenheng Tang, Nuo Chen, Wenxuan Wang, and Bingsheng He. Reasoning models can be easily hacked by fake reasoning bias. In <i>Lock-LLM Workshop: Prevent Unauthorized Knowledge Use from Large Language Models</i> .	960
		961
		962
		963
		964
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	965
		966
		967
		968
		969
		970
	Zhihui Xie, Jizhou Guo, Tong Yu, and Shuai Li. 2024. Calibrating reasoning in language models with internal consistency. <i>Advances in Neural Information Processing Systems</i> , 37:114872–114901.	971
		972
		973
		974
	Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-Li Lee, and Wynne Hsu. 2024. Faithful logical reasoning via symbolic chain-of-thought. <i>arXiv preprint arXiv:2405.18357</i> .	975
		976
		977
		978
	Nan Xu and Xuezhe Ma. 2024. Decoprompt: Decoding prompts reduces hallucinations when large language models meet false premises. <i>arXiv preprint arXiv:2411.07457</i> .	979
		980
		981
		982
	Yihao Xue, Kristjan Greenewald, Youssef Mroueh, and Baharan Mirzasoleiman. 2025. Verify when uncertain: Beyond self-consistency in black box hallucination detection. <i>arXiv preprint arXiv:2502.15845</i> .	983
		984
		985
		986
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025a. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	987
		988
		989
		990
		991
	Yijin Yang, Cristina Cornelio, Mario Leiva, and Paulo Shakarian. 2025b. Error detection and correction for interpretable mathematics in large language models. <i>arXiv preprint arXiv:2508.03500</i> .	992
		993
		994
		995

996	Zijun Yao, Yantao Liu, Yanxu Chen, Jianhui Chen, Junfeng Fang, Lei Hou, Juanzi Li, and Tat-Seng Chua. 2025. Are reasoning models more prone to hallucination? <i>arXiv preprint arXiv:2505.23646</i> .	1051
997		1052
998		1053
999		1054
1000	Weiqiu You, Anton Xue, Shreya Havaldar, Delip Rao, Helen Jin, Chris Callison-Burch, and Eric Wong. 2025. Probabilistic soundness guarantees in llm reasoning chains. <i>arXiv preprint arXiv:2507.12948</i> .	1055
1001		1056
1002		
1003		
1004	Hongbang Yuan, Pengfei Cao, Zhuoran Jin, Yubo Chen, Daojian Zeng, Kang Liu, and Jun Zhao. 2024. Whispers that shake foundations: Analyzing and mitigating false premise hallucinations in large language models. <i>arXiv preprint arXiv:2402.19103</i> .	1057
1005		1058
1006		1059
1007		1060
1008		1061
1009	Michael J Zellinger, Rex Liu, and Matt Thomson. 2025. Cost-saving llm cascades with early abstention. <i>arXiv preprint arXiv:2502.09054</i> .	1062
1010		1063
1011		1064
1012	Qingcheng Zeng, Weihao Xuan, Leyang Cui, and Rob Voigt. 2025. Thinking out loud: Do reasoning models know when they're right? <i>arXiv preprint arXiv:2504.06564</i> .	1065
1013		1066
1014		1067
1015		1068
1016	Anqi Zhang, Yulin Chen, Jane Pan, Chen Zhao, Aurojit Panda, Jinyang Li, and He He. 2025a. Reasoning models know when they're right: Probing hidden states for self-verification. In <i>Second Conference on Language Modeling</i> .	1069
1017		1070
1018		
1019		
1020		
1021	Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. 2024a. Rest-mcts*: Llm self-training via process reward guided tree search. <i>Advances in Neural Information Processing Systems</i> , 37:64735–64772.	1071
1022		1072
1023		1073
1024		1074
1025		1075
1026	Hanning Zhang, Shizhe Diao, Yong Lin, Yi Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng Ji, and Tong Zhang. 2024b. R-tuning: Instructing large language models to say 'I don't know'. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 7113–7139, Mexico City, Mexico.	1076
1027		1077
1028		1078
1029		1079
1030		1080
1031		1081
1032		1082
1033		
1034	Jiaxin Zhang, Zhuohang Li, Kamalika Das, Bradley Malin, and Sricharan Kumar. 2023a. Sac3: reliable hallucination detection in black-box language models via semantic-aware cross-check consistency. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 15445–15458.	1083
1035		1084
1036		1085
1037		1086
1038		1087
1039		1088
1040	Yongheng Zhang, Xu Liu, Ruihan Tao, Qiguang Chen, Hao Fei, Wanxiang Che, and Libo Qin. 2025b. Vitcot: Video-text interleaved chain-of-thought for boosting video understanding in large language models. In <i>Proceedings of the 33rd ACM International Conference on Multimedia</i> , pages 5267–5276.	1089
1041		
1042		
1043		
1044		
1045		
1046	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, and 1 others. 2025c. Siren's song in the ai ocean: A survey on hallucination in large language models. <i>Computational Linguistics</i> , pages 1–46.	1051
1047		1052
1048		1053
1049		1054
1050		1055
		1056
		1057
		1058
		1059
		1060
		1061
		1062
		1063
		1064
		1065
		1066
		1067
		1068
		1069
		1070
		1071
		1072
		1073
		1074
		1075
		1076
		1077
		1078
		1079
		1080
		1081
		1082
		1083
		1084
		1085
		1086
		1087
		1088
		1089