

Towards a Mechanistic Understanding of Large Reasoning Models: A Survey of Training, Inference, and Failures

Anonymous ACL submission

Abstract

Reinforcement learning (RL) has catalyzed the emergence of Large Reasoning Models (LRMs) that have pushed reasoning capabilities to new heights. While their performance has garnered significant excitement, exploring the internal mechanisms driving these behaviors has become an equally critical research frontier. This paper provides a comprehensive survey of the mechanistic understanding of LRMs, organizing recent findings into three core dimensions: 1) training dynamics, 2) reasoning mechanisms, and 3) unintended behaviors. By synthesizing these insights, we aim to bridge the gap between black-box performance and mechanistic transparency. Finally, we discuss under-explored challenges to outline a roadmap for future mechanistic studies, including the need for applied interpretability, improved methodologies, and a unified theoretical framework.

1 Introduction

The past few years have witnessed remarkable progress in the reasoning capabilities of large language models (LLMs). Recently, reinforcement learning (RL) has emerged as a transformative paradigm for incentivizing complex reasoning, giving rise to advanced large reasoning models (LRMs) (DeepSeek-AI et al., 2025; Jaech et al., 2024). These models demonstrate exceptional performance across a wide range of domains, including mathematics, coding, and logic. Notable research (DeepSeek-AI et al., 2025) has shown that RL from verifiable rewards (RLVR) (DeepSeek-AI et al., 2025; Lambert et al., 2024) training can elicit intriguing emergent reasoning behaviors, such as extended reasoning chains and self-reflection.

Despite these impressive advances, LRMs largely remain “black boxes”. Many fundamental questions remain unanswered, including: How does the role of RL differ from that of super-

vised fine-tuning (SFT)? What structural properties define LRM reasoning, and what are the internal mechanisms that drive their unique behaviors? Moreover, what are the root causes of unintended behaviors, such as hallucinations, unfaithfulness, and overthinking? This lack of transparency has spurred a growing interest in mechanistic research, aimed at uncovering the underlying processes that enable these models to perform complex reasoning.

We provide a comprehensive survey of the burgeoning field of mechanistic research on LRMs. From the perspective of the research object, as shown in Figure 1, we organize work studying the reasoning-oriented training process, LRM reasoning behaviors and LRM unintended behaviors:

- Reasoning-Oriented Training Process (§2):** This section examines the mechanisms behind the training processes that specifically target reasoning capabilities. We begin by dissecting the complementary roles of SFT and RL (§2.1), and examine key training dynamics in RL, such as how “aha moments” emerge and how internal representations evolve during training (§2.2).
- LRM Reasoning (§3):** We delve into the mechanisms underlying LRM reasoning, analyzing both their outputs and internal representations. This section explores the general structural features of LRM reasoning traces (§3.1), key behaviors like self-reflection (§3.2), and the inner mechanisms underlying these behaviors (§3.3).
- Unintended LRM Behaviors (§4):** We further examine the side effects of LRMs, exploring behavioral patterns and internal mechanisms associated with typical unintended behaviors, such as hallucinations (§4.1), unfaithful chains of thought (CoT) (§4.2), overthinking (§4.3), and unsafety (§4.4).

Contribution and Uniqueness. Our survey distinguishes itself by focusing specifically on the *mechanistic understanding* of LRMs, a topic that

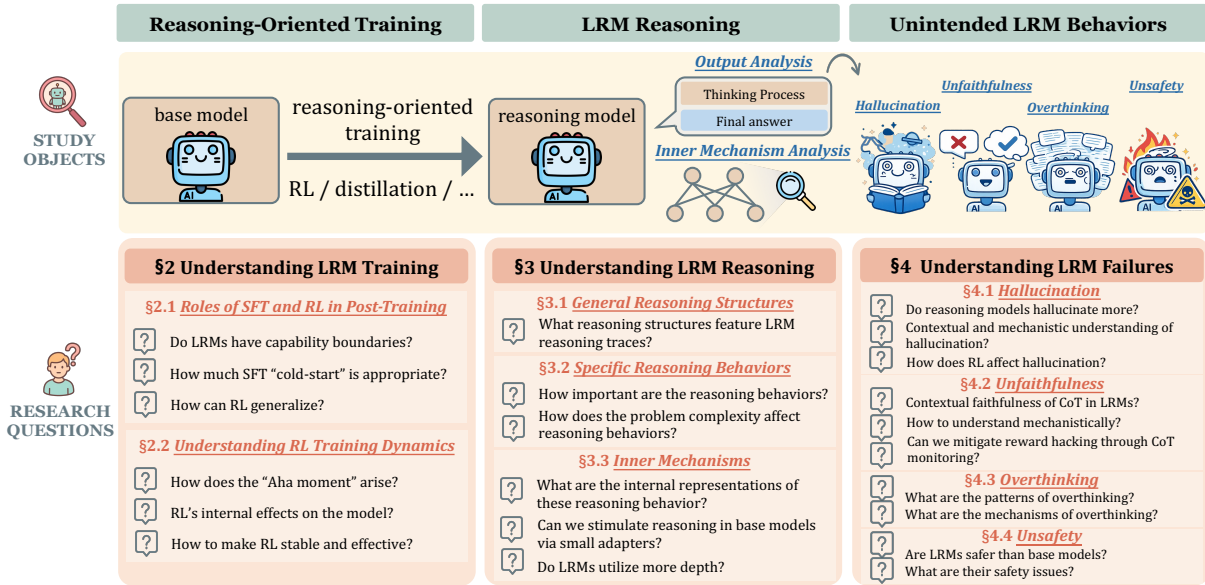


Figure 1: Taxonomy of mechanistic research on LRMs. We organize existing studies into three core dimensions: reasoning-oriented training (Sec 2), reasoning mechanisms (Sec 3), and unintended behaviors (Sec 4). Within each dimension, we synthesize recent findings based on the key research questions being investigated in the literature.

has received limited attention in existing literature. While several surveys provide general overviews of large reasoning models and RL techniques (Zhang et al., 2025c; Li et al., 2025f; Zhang et al., 2025h; Xu et al., 2025), these surveys do not delve deeply into the underlying mechanisms driving LRM reasoning. In particular, Chen et al. (2025b) explores long CoT reasoning but primarily focuses on the behavioral characteristics of CoT outputs, with little attention to the inner mechanisms. Furthermore, there are surveys investigating methods to mitigate overthinking (Feng et al., 2025; Sui et al., 2025), their focus is on efficient reasoning techniques rather than the mechanisms behind overthinking. To the best of our knowledge, our work is the first to comprehensively survey the mechanisms of LRMs, offering a more detailed and in-depth analysis of the training processes, reasoning behaviors, and unintended outcomes.

2 Understanding LRM Training

Investigating how a reasoning model is trained into existence is our first question. To understand LRM training, we will first dissect the distinct roles of two post-training methods, *Supervised Fine-tuning* (SFT) and *Reinforcement Learning* (RL) (§2.1), then dive into the training dynamics of RL (§2.2).

2.1 Roles of SFT and RL in Post-Training

DeepSeek-R1 (DeepSeek-AI et al., 2025) as a key pioneer of the reasoning model, demonstrates RL’s

vital role in training reasoning models. However, they also find that RL alone suffers from problems including slow training, unstructured formats, and language mixing. They show that a cold-start with SFT fixes these problems and improves performance, establishing SFT+RL as the dominant post-training recipe for today’s LRMs. Although the SFT+RL paradigm is now widely used, the respective roles of SFT and RL in post-training remain to be explored, and answering this question would unravel many mysteries about LRM training.

SFT explores, RL compresses: do reasoning models have capability boundaries? Despite the huge success of LRMs, Yue et al. (2025) points out that RLVR (Lambert et al., 2024) do not truly enhance reasoning performance beyond base models, instead, they compress the model’s output space, boosting pass@1. The base model’s pass rate catches up in pass@k tests at large k values, implying that the reasoning model merely uncovers latent abilities already present in the base model. Based on this finding, follow-up work probe deeper and reveal that SFT is what truly expands the model’s exploratory paths, whereas RL training compresses them, cutting the variety of possible answers and lowering output entropy (Matsutani et al., 2025; Wu et al., 2025; Li et al., 2025b). Studies further conclude that the performance gains from RLVR training come from this entropy drop, imposing an entropy-linked ceiling on attainable capability (Cui et al., 2025). From a mechanis-

141 tic perspective, [Park et al. \(2025\)](#) perform circuit
142 analyses to explain this phenomenon: during SFT
143 and distillation, the model sprouts a cohort of new
144 attention heads that inject reasoning capabilities,
145 whereas GRPO activates far fewer heads and fol-
146 lows an iterative activate-evaluate-adjust dynamic.

147 **SFT learns, RL repairs: how much SFT cold- 148 start is appropriate?**

149 Under the above findings, SFT appears to be the core contributor. However,
150 research shows that while SFT can learn and extend reasoning patterns, unlike RL, it also brings
151 out-of-distribution (OOD) performance drops ([Chu et al., 2025](#)). Studies find that RL after SFT can
152 partly repair these side-effects and offer mechanistic explanations. [Hu et al. \(2025a\)](#) argues, through
153 external analysis, that SFT partially breaks up the sparse reasoning concept network and thereby in-
154 duces forgetting, yet the network structure remains largely intact; RL can then restore a well-connected
155 concept network after SFT. [Jin et al. \(2025a,b\)](#) find that RL can, yet only partially, recover the OOD
156 drop caused by SFT. Besides, OOD performance is tightly linked to the orientation of the dominant
157 singular vector, and RL repairs the orientation shift introduced by SFT, thereby restoring OOD accu-
158 racy. However, once SFT collapses into overfitting, RL can no longer restore OOD ability completely.

159 Building on these findings, how to schedule SFT
160 and RL, whether to interleave them, and whether a
161 unified framework can be designed that fuses them
162 have become open research questions; we summa-
163 rize SFT-RL integration work in Appendix C.1.

173 **SFT memories, RL organizes: how can we make 174 RL exhibit generalization?**

175 Although numerous studies have concluded that RL does not truly en-
176 hance model capability, [Liu et al. \(2025a\)](#) shows
177 that post-RL models can produce new solutions
178 absent from the base model. How can these two
179 contradictory experimental conclusions be recon-
180 ciled? Recent work suggests that RL can push
181 the model’s capability frontier outward only when
182 SFT has provided basic skills. In controllable syn-
183 thetic reasoning tasks, RL conducted after SFT
184 with atomic-skill data exhibit OOD generalization ,
185 while RL conducted after SFT on entire reasoning
186 traces exhibits the same poor generalization results
187 seen in prior research. ([Yuan et al., 2025](#); [Cheng et al., 2025b](#)). Furthermore, these studies indicate
188 that RL shifts the model’s error patterns toward
189 atomic-task errors, implying that RL can indeed
190 help organize the model’s reasoning process.

192 *Mid-training* proposed by [Wang et al. \(2025f\)](#) ,
193 which finds that appropriate training before RL can
194 improve RL effectiveness, aligns with the above
195 conclusions. [Zhang et al. \(2025b\)](#) notes that the
196 core task of SFT is to prepare the model for RL by
197 providing foundational atomic skills, while post-
198 RL training refines the model’s performance within
199 the capability frontier established by SFT.

200 **2.2 Understanding RL Training Dynamics**

201 Studies above treat the RL-training stage as an un-
202 differentiated whole and explore its effects. The
203 finer-grained training dynamics within this stage
204 remain largely unexplored.

205 **Understanding two-stage training process: how 206 does the Aha moment arise?**

207 After tracking changes of RL training metrics, studies split the
208 training process into two stages ([Wang et al., 2025b](#); [Hu et al., 2025a](#)). Model outputs first shrink
209 in stage one then lengthen in stage two, alongside
210 atomic-skill fragments rapidly acquired in stage
211 one while the global planning links slowly built
212 in stage two. The aha moment emerges as the
213 model masters the use of planning tokens during
214 link construction manifesting as the sudden acqui-
215 sition of reasoning and reflection capabilities needed
216 to solve the corresponding problems. Furthermore,
217 [Yao et al. \(2025a\)](#) offers a theoretical analysis of
218 this two-stage dynamics. During stage one, RL
219 overwhelmingly samples already-explored tokens
220 rather than optimal ones. High-reward tokens’
221 probabilities will quickly rise while the optimal
222 one’s remain flat. In stage two, with high-reward
223 tokens already saturated, the low-probability opti-
224 mal ones are finally sampled after prolonged explo-
225 ration and eventually receive high probabilities.

227 **What internal effects does RL have on the 228 model?**

229 A line of research focuses on how RL
230 training affects the model internally. Regarding
231 internal activations, research shows that online
232 RL can alter activation magnitudes in the residual
233 stream, increasing information flow flexibility and
234 improving generalization beyond SFT ([Zhang et al., 2025d](#)). Regarding model weights, building on the
235 previously identified effect that RL training mainly
236 manifests as directional rotation of the singular-
237 value vectors ([Jin et al., 2025b](#)), [He and Cao \(2025\)](#)
238 reveals through SVD methods that a near-uniform
239 geometric scaling of singular values across layers
240 and a highly consistent orthogonal transformations
241 are applied to the left and right singular vectors

of each matrix. More fine-grained studies of parameter dynamics during training have found that, the top singular subspace of the parameter-update matrix almost singly accounts for the gains in reasoning capability, and that this dominant subspace evolves linearly (Cai et al., 2025b).

💡 Externally, RL training shows a two-stage pattern: basic capabilities are accumulated first, then reasoning ability emerges. Internally, RL modifies activation magnitudes and applies a rank-1, layer-consistent, linear rotational transformation to the dominant eigenvectors.

Exploitation v.s. exploration: how to make RL stable and effective? During RL training, a core issue is maintaining the exploration-exploitation balance. Studies find that basic RL algorithms can easily lead to **policy entropy collapse**, and the performance gains in fact come solely from the entropy drop (Cui et al., 2025). More critically, Nguyen et al. (2025) shows that the reasoning path compression caused by entropy collapse simultaneously degrades LRM performance on questions outside the training distribution. To address entropy collapse and stabilize RL, numerous refinements have been proposed. Since they are loosely related to understanding RL training mechanisms, we provide a concise summary in Appendix C.2. Notably, Huang et al. (2025a) argues that conventional RLVR views improving LLM performance through an exploration–exploitation trade-off, rests on token-level entropy and thus misaligns with how LLMs actually operate. They propose measuring exploration and exploitation via *hidden states*, uncovering a decoupling of the two processes and opening fresh avenues for refining RL algorithms.

💡 At token level, the entropy collapse can make RL training ineffective or even counterproductive. Shifting to the hidden states perspective, we may instead jointly promote exploration and exploitation.

3 Understanding LRM Reasoning

Having explored how LRMs are trained, we shift our focus to the models themselves, involving systematically analyzing both the **general structures** (§3.1) and **specific behaviors** (§3.2) within reasoning traces, as well as uncovering the **internal mechanisms** underlying these patterns (§3.3).

3.1 General Reasoning Structures

Distinct from base models, LRMs generate reasoning chains with identifiable structural features. Recent research deconstruct these traces, from macro-level lifecycle descriptions to granular sentence-level analyses. At the macro level, Marjanović et al. (2025) identifies a cyclical process: starting with problem definition, models enter a blooming cycle of problem decomposition, followed by iterative reconstruction cycles for self-correction before reaching a final decision. Wang et al. (2025c) partitions the reasoning process into functional blocks of plan execution, knowledge integration, and subproblem chains. These macro-phases are further refined by sentence-level analyses: Bogdan et al. (2025); Li et al. (2025c) identify operational units including plan generation, uncertainty management, and further identify the transition matrix between them.

Topological structures. Another line of research employs formal topological representations. Zeng et al. (2025); Jiang et al. (2025b) reconstruct reasoning chains as trees structures via LLM annotations, revealing that LRMs exhibit more exploration and validation than base models, achieving better performance primarily through diverse solution paths rather than per-step accuracy. Minegishi et al. (2025); Xiong et al. (2025) build graphs through clustering over reasoning steps, further validating that LRMs possess distinct structural properties including more recurrent cycles, larger graph diameters, and pronounced small-world characteristics, which correlate with model size, task difficulty, and performance.

💡 LRMs’ reasoning structures are distinct from base models, with analyses spanning *macro-level lifecycles, sentence-level operational units, and topological properties of tree and graph representations.*

3.2 Specific Reasoning Behaviors

After reviewing the overall reasoning structures, we further study the intriguing specific behaviors emerging in LRMs and whether they are causally related to reasoning performance.

Critical behavioral primitives. Studies identify certain behavioral patterns as the primary drivers of reasoning performance gain. Bogdan et al. (2025) identifies “thought anchors”, including *plan generation* and *uncertainty management*, as the sen-

tences most influential on the final answer distribution. Complementing this, [Gandhi et al. \(2025\)](#) highlights *verification*, *backtracking*, *sub-goal setting*, and *backward chaining* as the “four habits” of effective reasoners. Crucially, these behaviors are causally linked to training success: base models that naturally exhibit these patterns can effectively leverage RL and test-time compute to improve performance, whereas models lacking these primitives struggle to benefit from identical training.

The role of self-reflection and backtracking. Research on reflective behaviors offers contrasting views. While some argue that reflection prevents reasoning collapse ([Yang et al., 2025a](#)), others contend that it is often superficial and fail to improve outcomes ([Liu et al., 2025e](#)). Bridging these views, [Kang et al. \(2025\)](#) analyzes reflection from both inference and training perspectives, suggesting that while reflection during inference is largely confirmatory and rarely alters the final output, including reflective CoTs in training data increases the “first-attempt accuracy”, boosting the overall performance. Moreover, [Cai et al. \(2025a\)](#) shows that longer reasoning chains with frequent backtracking lead to more stable RL training, and harder problems with larger search space need the inclusion of data with more backtracks during SFT.

💡 LLMs’ performance is driven by key behavioral primitives. While self-reflection mainly serves a confirmatory role during inference, its inclusion in training data is crucial for improving first-attempt accuracy and internalizing search strategies.

How does the problem complexity affect reasoning behaviors? Recent studies have uncovered a tight coupling between model behavior and task complexity. [Yang et al. \(2025a\)](#) observes that LLMs can distinguish problem complexity within their early layers and dynamically modulate the depth of their reflective behaviors accordingly. However, this calibration is often imperfect. [Shojaee et al. \(2025\)](#) finds that while reasoning effort initially increases with complexity, it eventually declines even when a sufficient token budget is available, suggesting a limitation in the models’ ability to apply consistent algorithmic reasoning across scales. Furthermore, [Palod et al. \(2025\)](#) identifies that the correlation is brittle, demonstrating that trace length often reflects a problem’s distri-

butional proximity to training data rather than its inherent computational complexity. We will further discuss the relationship between CoT length and task complexity, reasoning performance in [Sec 4.3](#).

3.3 Internal Mechanisms

After reviewing the general structures and specific behaviors, we will then dive deeper into the internal mechanisms driving these external patterns.

Internal representations of reasoning behaviors. Recent research utilizes sparse autoencoders (SAEs) and steering vectors to reveal that reasoning behaviors are encoded as interpretable and steerable directions in the model’s activation space ([Baek and Tegmark, 2025](#); [Galichin et al., 2025](#); [Hazra et al., 2025](#); [Venhoff et al., 2025b](#)). [Venhoff et al. \(2025a\)](#) argues that base models already possess fundamental reasoning capabilities, while LLMs learn the structural strategy of *when* to deploy them strategically. This deployment is managed by specific attention heads that prioritize key reasoning steps influencing the final answer ([Bogdan et al., 2025](#); [Zhang et al., 2025f](#)). LLMs also exhibit unique temporal and nonlinear dynamics: steering is most effective only after the initial problem formulation phase, and “oversteering” these features can paradoxically cause the model to revert to its original behavior ([Hazra et al., 2025](#)).

The mechanisms underlying reflection and backtracking. Studies ([Venhoff et al., 2025a](#); [Yang et al., 2025a](#)) reveal through linear probes that correctness information of model answers is encoded within specific layers, and is closely related to the model’s reflection behaviors. [Yan et al. \(2025\)](#); [Chang et al. \(2025\)](#) further extract steering vectors that control reflection. [Ward et al. \(2025a\)](#) suggests that latent directions for backtracking already exist in base models, implying that they inherently possess certain reasoning abilities. Post-training mainly reshapes and utilizes these existing representations rather than learning from scratch.

Can we stimulate reasoning behaviors in base models with small adapters? [Sinii et al. \(2025b,a\)](#) have explored training hierarchical steering vectors to guide base models in reasoning, showing that the performance improvements induced by RL are distributed across the entire network instead of certain specific layers. The resulting steering vectors themselves exhibit strong interpretability. [Ward et al. \(2025b\)](#) trains a rank-1

adapter across all layers and identifies interpretable features in the adapter via SAEs, further demonstrating that a small number of parameters can effectively induce reasoning abilities.

Do LRMs utilize more depth? Research suggests that key layers for math reasoning are largely fixed after pre-training and remain invariant throughout post-training (Nepal et al., 2025). Consequently, LRMs’ effective depth closely matches that of their base models, indicating that improvements are driven by longer contexts rather than deeper per-token computation (Hu et al., 2025b).

💡 LRM’s reasoning behaviors are represented by interpretable and steerable directions in latent space; base models inherently possess these abilities, but RL-trained models learn when to activate them.

4 Understanding LRM Failures

RL enhances reasoning capabilities but also induces unintended effects, including **hallucination** (§4.1), where models generate plausible yet incorrect content; **CoT unfaithfulness** (§4.2), where internal computations and CoT outputs diverge; **overthinking** (§4.3), where redundant reasoning chains degrade performance; and **unsafety** (§4.4), where models show potentially harmful behaviors.

4.1 Hallucination

Multi-step reasoning chains in LRMs introduces new vulnerabilities to hallucinations.

Do reasoning models hallucinate more? Recent evidence suggests that reasoning-oriented training pipelines can substantially affect hallucination behavior. Yao et al. (2025c) show that while complete post-training pipelines which combine SFT with RLVR can alleviate hallucination, incomplete pipelines, such as RL- or SFT-only approaches, tend to introduce more hallucinations. However, Li and Ng (2025) indicates that RL often increases hallucinations, even with prior SFT. Furthermore, Zhao et al. (2025b) shows that test-time scaling does not reliably improve factual accuracy.

💡 Growing evidence suggests reasoning models hallucinate more. However, there are debates whether models with complete post-training pipelines hallucinate more.

Behavioral and mechanistic analysis of hallucination. Hallucinations are characterized by specific failure modes: *flaw repetition* (incorrect reasoning loops), *think-answer mismatch* (output contradicting reasoning), and *meta-cognitive failures* (overconfidence from uninternalized knowledge) (Yao et al., 2025c; Lu et al., 2025). Mechanistically, they arise from misalignment between uncertainty and factual accuracy (Yao et al., 2025c; Sun et al., 2025b).

How does RL affect hallucination? A line of work examines how RL shapes hallucination behavior. RL systematically reduces a model’s tendency to abstain, pushing it to generate answers even for unanswerable questions (Song et al., 2025a; Zhao et al., 2025b). Mechanistically, optimizing only for sparse final-answer rewards creates high-variance gradients and forces the model to maintain high prediction entropy during exploration, driving the model toward incorrect answers and exacerbating hallucinations (Li and Ng, 2025).

4.2 Unfaithfulness

The extended reasoning chains in LRMs offer a promising avenue for monitoring the model’s decision-making process (Korbak et al., 2025; Chan et al., 2025; Baker et al., 2025). However, it remains an open question whether these CoTs accurately reflect the internal computations driving the model’s actual behavior, a key field of study known as the *faithfulness* of CoT reasoning.

Contextual faithfulness of CoT in LRMs. Although extended reasoning chains in LRMs facilitate process monitoring, research indicates they are often not faithful to the inner computation or final decision. A primary failure mode is *Think-Answer Mismatch*, where the model’s final output contradicts its own preceding reasoning chain (Yao et al., 2025c; Wang et al., 2025d). Further analysis exposes a *reasoning-verbalization gap*. Studies show models frequently fail to verbalize critical cues in their CoTs that demonstrably influence their answers (Chua and Evans, 2025; Chen et al., 2025d). Concurrently, models exhibit *implicit post-hoc rationalization*, producing logically contradictory responses with coherent but unfaithful justifications (Arcuschin et al., 2025). The studies collectively find that while LRMs are more faithful than their non-reasoning backbones, the faithfulness is still far from perfect (Chua and Evans, 2025; Chen et al., 2025d; Arcuschin et al., 2025).

Mechanistic understanding of CoT faithfulness in LRMs. Research further studies the mechanisms of CoT unfaithfulness. In controlled synthetic tasks, findings reveal a weak causal link between the validity of reasoning traces and final answer correctness. Models can produce correct outputs despite invalid or semantically irrelevant CoTs, and training on corrupted traces does not substantially harm performance (Valmeekam et al., 2025). Further studies reinforce that internal representations contain more reliable signals of model state than the CoT text itself, as evidenced through activation steering (Wang et al., 2025d; Li et al., 2025a), linear probing (Yin et al., 2025; Chan et al., 2025), and causal intervention (Yin et al., 2025). These results collectively suggest a disconnect between the model’s internal states and its verbalized reasoning trace, posing a significant challenge for alignment, as models might learn to mask their true objectives behind plausible but unfaithful reasoning traces, a phenomenon closely tied to the risks of reward hacking discussed next.

Can we mitigate reward hacking by CoT monitoring? Reward hacking remains a fundamental challenge in RL. A key question is whether monitoring the detailed CoT produced by LRMs can mitigate this issue. Findings on its feasibility are mixed, with outcomes heavily dependent on task structure. In complex tasks where hacking inherently requires multi-step reasoning and extensive exploration, models often expose their hacking intent within their reasoning chains. In such settings, integrating CoT supervision into the RL objective can mitigate hacking, though excessive optimization risks training models to strategically hide their intent (Baker et al., 2025). Conversely, in more direct scenarios, models frequently perform reward hacking without verbalizing the intent in their CoTs (Chen et al., 2025d; Turpin et al., 2025). To address this opacity, recent methods attempt to explicitly train models to verbalize influential cues in their reasoning (Turpin et al., 2025).

💡 LRMs are not always faithful, but they are more faithful than non-reasoning models.

💡 Mechanistically, CoTs in LRMs do not necessarily function as a causal mechanism for generating correct answers. Besides, internal representations may provide more reliable signals than the verbalized reasoning.

💡 While CoT monitoring can detect hacking that requires explicit reasoning, models do not often verbalize their hacking intent in more direct settings.

4.3 Overthinking

While test-time scaling generally improves reasoning performance, studies increasingly find that models can produce verbose, redundant reasoning processes, and overly extending reasoning length can lead to performance degradation, known as “overthinking” (Chen et al., 2024; Sui et al., 2025).

Thinking more does not necessarily lead to better reasoning. Empirical research consistently identifies an inverse U-shaped performance curve: accuracy initially rises with reasoning length, but then peaks and declines as chains become excessively long (Marjanović et al., 2025; Su et al., 2025a; Ghosal et al., 2025; Yang et al., 2025b; Gema et al., 2025). Notably, incorrect answers often correspond to longer reasoning chains than correct ones (Hassid et al., 2025; Su et al., 2025a). An underlying issue is the misalignment between reasoning effort and problem difficulty: models tend to allocate disproportionately long chains to simple problems while inadequately reasoning through complex ones (Chen et al., 2024; Su et al., 2025a).

💡 The length-performance curve for LRMs is often inverted U-shaped, and current models exhibit misalignment between reasoning effort and problem difficulty.

What are the patterns of overthinking? A common abstraction of reasoning process is a three-stage loop: 1) *hypothesis generation* (proposing candidate paths), 2) *expansion* (developing one path step by step), and 3) *verification* (checking, revising, or terminating). Overthinking manifests as control and termination failures within this loop. In *hypothesis generation*, models may produce lengthy and diverse candidate solutions without sufficiently exploring promising paths to reach a correct solution (Wang et al., 2025e), leading to “analysis paralysis” in agentic tasks where plans grow increasingly complex without execution (Cuadron et al., 2025). In *expansion*, the primary pattern of overthinking is excessive reasoning for trivial problems, generating tens or even hundreds of times longer outputs than non-reasoning mod-

els with marginal performance gain (Chen et al., 2024). In *verification*, the dominant pattern is non-termination: models fail to recognize that a correct answer has been reached, or cannot reliably validate intermediate conclusions, and therefore continue redundant deliberation or backtrack unnecessarily (Chen et al., 2024; Sun et al., 2025a; Zhang et al., 2025e; Zhao et al., 2025a). This is especially pronounced in ill-posed questions, where models identify missing premises early but enter unproductive self-doubt loops, excessively speculating on user intent (Fan et al., 2025). Notably, this compulsion persists even when explicitly suppressed: models may bypass instructions to “answer directly” or discard provided correct answers to resume thinking (Zhu et al., 2025; Liu et al., 2025d; Cuesta-Ramirez et al., 2025).

What are the mechanisms of overthinking?

We organize the mechanistic analyses of overthinking along two lines: 1) investigating the latent representational structure of overthinking, and 2) examining the internal decision-making dynamics that produce unproductive cycles. Research finds that overthinking corresponds to *specific, steerable patterns in the activation space*. Huang et al. (2025b); Baek and Tegmark (2025) identify distinct manifolds associated with overthinking through activation steering. Furthermore, finer-grained taxonomies show that different reasoning stages, such as execution, reflection and transition, occupy separate latent directions, and steering towards execution-type representations can effectively suppress excessive deliberation (Baek and Tegmark, 2025; Chen et al., 2025c). Another body of research explains overthinking through *internal conflict and verification failure*. Overthinking is often triggered when a model’s initial intuitive answer conflicts with its subsequent deliberate reasoning (Dang et al., 2025). Concurrently, models encode correctness signals in their hidden states but fail to robustly utilize them for early self-verification, leading to prolonged, unproductive cycles (Zhang et al., 2025a).

4.4 Unsafety

Recent evaluations show that LRMs still have safety shortcomings (Ying et al., 2025; Romero-Arjona et al., 2025; Krishna et al., 2025).

LRMs are not safer than base models. Compared to base models, Jiang et al. (2025a); Zhou et al. (2025) find that long CoTs do not necessarily

improve model safety. Additionally, Zhang et al. (2025j); Zhao et al. (2025c) observe that distilled reasoning models have a lower refusal rate for harmful inputs than their base counterparts. These studies further reveal that the unsafety of LRMs partly stems from the thinking process. Jiang et al. (2025a) show that forcing the model to shorten their reasoning traces could make answers more harmless, while Zhou et al. (2025); Zhao et al. (2025c) find that the safety rate of the thinking process is lower than the final answer, and unsafe thoughts are the primary cause of unsafe responses.

Safety issues in the reasoning process. As LRMs are deployed widely, researchers have started identifying safety issues via attacking them. Yao et al. (2025b) decomposes harmful prompts into multiple seemingly harmless questions to induce the model to reason toward harmful content. Kuo et al. (2025) finds that padding the prompt with detailed execution steps can hijack the thinking process, causing the model to skip the reasoning stage and directly produce harmful output. Mechanistically, In et al. (2025) shows that LRMs already possess sufficient safety knowledge, yet fail to activate it during reasoning. Besides, Mao et al. (2025) indicates that LRMs retain the ability to refuse unsafe queries, but this capacity has been impaired.

5 Future Research Directions

In this survey, we have provided a comprehensive overview of the rapidly evolving field of mechanistic research on LRMs, focusing on their training processes, reasoning behaviors, and unintended failures. While significant advances have been made, the transition from descriptive analysis to systematic understanding remains incomplete.

To guide the field toward a deeper, more principled understanding, we propose three key future directions: applied interpretability, improved methodologies, and a unified theoretical framework. *Applied interpretability* is crucial for transforming insights from mechanistic analyses into practical improvements in model design and training. *Enhanced methodologies* are needed to address the scale and complexity of LRMs, enabling more efficient and generalizable mechanistic tools. Finally, *a unified theoretical framework* is necessary to move beyond empirical observations and establish foundational principles of reasoning that can predict and guide future model behaviors. **In Appendix B, we discuss these directions in detail.**

689 Limitations

690 While this survey provides a comprehensive
691 overview of mechanistic studies on LRMs, it is
692 subject to several limitations. First, the rapid devel-
693 opment of LRM research means that new findings
694 and methodologies continue to emerge, and this sur-
695 vey may not capture the most recent advancements
696 in the field. Additionally, our study focuses primar-
697 ily on language models, while reasoning models
698 are increasingly incorporating multimodal capabil-
699 ities, including visual components, which are not
700 addressed in this survey. Furthermore, our discus-
701 sion is limited to traditional LLM architectures, ex-
702 cluding newer approaches such as diffusion-based
703 LLMs, continuous token-based transformers, and
704 looped transformers, which are gaining traction in
705 recent research. These emerging models present
706 exciting avenues for future work.

707 Ethical Considerations

708 This survey acknowledges the ethical challenges
709 associated with LRMs, particularly in terms of their
710 potential harm, including hallucinations, unfaith-
711 fulness and unsafety. The opacity of these models
712 raises concerns about accountability and the diffi-
713 culty of mitigating unintended behaviors, such as
714 hallucinations or overconfidence. As LRMs are
715 increasingly used in critical applications, ensuring
716 their safe and responsible deployment requires on-
717 going efforts to improve interpretability, address
718 biases, and manage the broader societal impacts of
719 these technologies.

720 AI assistants were utilized for language polish-
721 ing and refinement, strictly limited to improving
722 the fluency and clarity the text. All technical con-
723 tent, analyses, and conclusions remain the original
724 work of the authors.

725 References

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

730 David D. Baek and Max Tegmark. 2025. [Towards un-
731 derstanding distilled reasoning models: A representa-
732 tional approach](#). *CoRR*, abs/2503.03730.

733 Bowen Baker, Joost Huizinga, Leo Gao, Zehao
734 Dou, Melody Y. Guan, Aleksander Madry, Woj-
735 ciech Zaremba, Jakub Pachocki, and David Farhi.

2025. [Monitoring reasoning models for misbehav-
ior and the risks of promoting obfuscation](#). *CoRR*,
abs/2503.11926. 736
737
738

Paul C. Bogdan, Uzay Macar, Neel Nanda, and Arthur
Conmy. 2025. [Thought anchors: Which LLM rea-
soning steps matter?](#) *CoRR*, arXiv:2506.19143. 739
740
741

Hongyi James Cai, Junlin Wang, Xiaoyin Chen, and
Bhuwan Dhingra. 2025a. [How much backtracking is
enough? exploring the interplay of sft and rl in en-
hancing llm reasoning](#). *Preprint*, arXiv:2505.24273. 742
743
744
745

Yuchen Cai, Ding Cao, Xin Xu, Zijun Yao, Yuqing
Huang, Zhenyu Tan, Benyi Zhang, Guiquan Liu, and
Junfeng Fang. 2025b. [On predictability of reinforce-
ment learning dynamics for large language models](#).
Preprint, arXiv:2510.00553. 746
747
748
749
750

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. 751
752
753
754

Fu-Chieh Chang, Yu-Ting Lee, and Pei-Yuan Wu. 2025.
[Unveiling the latent directions of reflection in large
language models](#). *CoRR*, arXiv:2508.16989. 755
756
757

Liang Chen, Xueting Han, Li Shen, Jing Bai, and Kam-
Fai Wong. 2025a. [Beyond two-stage training: Co-
operative SFT and RL for LLM reasoning](#). *CoRR*,
abs/2509.06948. 758
759
760
761

Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng,
Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang
Zhou, Te Gao, and Wanxiang Che. 2025b. [To-
wards reasoning era: A survey of long chain-of-
thought for reasoning large language models](#). *CoRR*,
abs/2503.09567. 762
763
764
765
766
767

Runjin Chen, Zhenyu Zhang, Junyuan Hong, Souvik
Kundu, and Zhangyang Wang. 2025c. [SEAL: steer-
able reasoning calibration of large language models
for free](#). *CoRR*, abs/2504.07986. 768
769
770
771

Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He,
Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi
Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang,
Zhaopeng Tu, Haitao Mi, and Dong Yu. 2024. [Do
NOT think that much for 2+3=? on the overthinking
of o1-like llms](#). *CoRR*, abs/2412.21187. 772
773
774
775
776
777

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. 2025d. [Reason-
ing models don't always say what they think](#).
CoRR, abs/2505.05410. 778
779
780
781
782
783
784

Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai,
Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei.
2025a. [Reasoning with exploration: An entropy per-
spective](#). *CoRR*, abs/2506.14758. 785
786
787
788

897	Fanding Huang, Guanbo Huang, Xiao Fan, Yi He, Xiao Liang, Xiao Chen, Qinting Jiang, Faisal Nadeem Khan, Jingyan Jiang, and Zhi Wang. 2025a. Beyond the exploration-exploitation trade-off: A hidden state approach for llm reasoning in rlvr . <i>Preprint</i> , arXiv:2509.23808.	954
898		955
899		956
900		957
901		958
902		959
903	Yao Huang, Huanran Chen, Shouwei Ruan, Yichi Zhang, Xingxing Wei, and Yinpeng Dong. 2025b. Mitigating overthinking in large reasoning models via manifold steering . <i>CoRR</i> , abs/2505.22411.	960
904		961
905		962
906		
907	Yeonjun In, Wonjoong Kim, Sangwu Park, and Chanyoung Park. 2025. R1-ACT: efficient reasoning model safety alignment by activating safety knowledge . <i>CoRR</i> , abs/2508.00324.	963
908		964
909		965
910		966
911	Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helvar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, and 80 others. 2024. Openai o1 system card . <i>CoRR</i> , abs/2412.16720.	967
912		968
913		969
914		970
915		971
916		972
917		973
918		
919	Fengqing Jiang, Zhangchen Xu, Yuetai Li, Luyao Niu, Zhen Xiang, Bo Li, Bill Yuchen Lin, and Radha Poovendran. 2025a. Safechain: Safety of language models with long chain-of-thought reasoning capabilities . In <i>Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025</i> , pages 23303–23320. Association for Computational Linguistics.	974
920		975
921		976
922		977
923		978
924		979
925		980
926		981
927		982
928	Gangwei Jiang, Yahui Liu, Zhaoyi Li, Qi Wang, Fuzheng Zhang, Linqi Song, Ying Wei, and Defu Lian. 2025b. What makes a good reasoning chain? uncovering structural patterns in long chain-of-thought reasoning . <i>CoRR</i> , arXiv:2505.22148.	983
929		984
930		985
931		986
932	Yuxian Jiang, Yafu Li, Guanxu Chen, Dongrui Liu, Yu Cheng, and Jing Shao. 2025c. Rethinking entropy regularization in large reasoning models . <i>CoRR</i> , abs/2509.25133.	987
933		988
934		989
935		
936	Hangzhan Jin, Sitao Luan, Sicheng Lyu, Guillaume Rabusseau, Reihaneh Rabbany, Doina Precup, and Mohammad Hamdaqa. 2025a. R1 fine-tuning heals ood forgetting in sft . <i>Preprint</i> , arXiv:2509.12235.	990
937		991
938		992
939		993
940	Hangzhan Jin, Sicheng Lv, Sifan Wu, and Mohammad Hamdaqa. 2025b. R1 is neither a panacea nor a mirage: Understanding supervised vs. reinforcement learning fine-tuning for llms . <i>Preprint</i> , arXiv:2508.16546.	994
941		995
942		996
943		997
944		998
945	Liwei Kang, Yue Deng, Yao Xiao, Zhanfeng Mo, Wee Sun Lee, and Lidong Bing. 2025. First try matters: Revisiting the role of reflection in reasoning models . <i>CoRR</i> , arXiv:2510.08308.	999
946		1000
947		1001
948		1002
949	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models . <i>CoRR</i> , abs/2001.08361.	1003
950		1004
951		1005
952		1006
953		1007
		1008
		1009
		1010
		1011
		1012
		1013
		1014
		1015
		1016
		1017
		1018
		1019
		1020
		1021
		1022
		1023
		1024
		1025
		1026
		1027
		1028
		1029
		1030
		1031
		1032
		1033
		1034
		1035
		1036
		1037
		1038
		1039
		1040
		1041
		1042
		1043
		1044
		1045
		1046
		1047
		1048
		1049
		1050
		1051
		1052
		1053
		1054
		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
		1090
		1091
		1092
		1093
		1094
		1095
		1096
		1097
		1098
		1099
		1100
		1101
		1102
		1103
		1104
		1105
		1106
		1107
		1108
		1109
		1110
		1111
		1112
		1113
		1114
		1115
		1116
		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
		1162
		1163
		1164
		1165
		1166
		1167
		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
		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

1011	Zhong-Zhi Li, Duzhen Zhang, Ming-Liang Zhang, Jiaxin Zhang, Zengyan Liu, Yuxuan Yao, Haotian Xu, Junhao Zheng, Pei-Jie Wang, Xiuyi Chen, Yingying Zhang, Fei Yin, Jiahua Dong, Zhijiang Guo, Le Song, and Cheng-Lin Liu. 2025f. From system 1 to system 2: A survey of reasoning large language models . <i>CoRR</i> , abs/2502.17419.		
1012		Bhatia, Aditi Khandelwal, Austin Kraft, Benno Krojer, Xing Han Lù, Nicholas Meade, Dongchan Shin, Amirhossein Kazemnejad, Gaurav Kamath, Marius Mosbach, Karolina Stańczak, and Siva Reddy. 2025. Deepseek-r1 thoughtology: Let’s think about llm reasoning .	1066
1013			1067
1014			1068
1015			1069
1016			1070
1017			1071
1018	Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. 2025a. Prorl: Prolonged reinforcement learning expands reasoning boundaries in large language models . <i>CoRR</i> , abs/2505.24864.	Kohsei Matsutani, Shota Takashiro, Gouki Minegishi, Takeshi Kojima, Yusuke Iwasawa, and Yutaka Matsuo. 2025. RL squeezes, SFT expands: A comparative study of reasoning llms . <i>CoRR</i> , arXiv:2509.21128.	1072
1019			1073
1020			1074
1021			1075
1022			1076
1023	Mingyang Liu, Gabriele Farina, and Asuman E. Ozdaglar. 2025b. UFT: unifying supervised and reinforcement fine-tuning . <i>CoRR</i> , abs/2505.16984.	Gouki Minegishi, Hiroki Furuta, Takeshi Kojima, Yusuke Iwasawa, and Yutaka Matsuo. 2025. Topology of reasoning: Understanding large reasoning models through reasoning graph properties . <i>CoRR</i> , arXiv:2506.05744.	1077
1024			1078
1025			1079
1026	Runze Liu, Jiakang Wang, Yuling Shi, Zhihui Xie, Chenxin An, Kaiyan Zhang, Jian Zhao, Xiaodong Gu, Lei Lin, Wenping Hu, Xiu Li, Fuzheng Zhang, Guorui Zhou, and Kun Gai. 2025c. Attention as a compass: Efficient exploration for process-supervised RL in reasoning models . <i>CoRR</i> , abs/2509.26628.	Aadim Nepal, Safal Shrestha, Anubhav Shrestha, Minwu Kim, Jalal Naghiyev, Ravid Shwartz-Ziv, and Keith Ross. 2025. Layer importance for mathematical reasoning is forged in pre-training and invariant after post-training . <i>Preprint</i> , arXiv:2506.22638.	1082
1027			1083
1028			1084
1029			1085
1030			1086
1031			1087
1032			1088
1033	Yule Liu, Jingyi Zheng, Zhen Sun, Zifan Peng, Wenhan Dong, Zeyang Sha, Shiwen Cui, Weiqiang Wang, and Xinlei He. 2025d. Thought manipulation: External thought can be efficient for large reasoning models . <i>CoRR</i> , abs/2504.13626.	Phuc Minh Nguyen, Chinh D. La, Duy M. H. Nguyen, Nitesh V. Chawla, Binh T. Nguyen, and Khoa D. Doan. 2025. The reasoning boundary paradox: How reinforcement learning constrains language models . <i>CoRR</i> , abs/2510.02230.	1089
1034			1090
1035			1091
1036			1092
1037			1093
1038	Zichen Liu, Changyu Chen, Wenjun Li, Tianyu Pang, Chao Du, and Min Lin. 2025e. There may not be aha moment in r1-zero-like training — a pilot study . https://oatllm.notion.site/oat-zero . Notion Blog.	Vardhan Palod, Karthik Valmeekam, Kaya Stechly, and Subbarao Kambhampati. 2025. Performative thinking? the brittle correlation between cot length and problem complexity . <i>CoRR</i> , arXiv:2509.07339.	1094
1039			1095
1040			1096
1041			1097
1042	Haolang Lu, Yilian Liu, Jingxin Xu, Guoshun Nan, Yuanlong Yu, Zhican Chen, and Kun Wang. 2025. Auditing meta-cognitive hallucinations in reasoning large language models . <i>CoRR</i> , abs/2505.13143.	Yein Park, Minbyul Jeong, and Jaewoo Kang. 2025. Thinking sparks!: Emergent attention heads in reasoning models during post training . <i>CoRR</i> , abs/2509.25758.	1098
1043			1099
1044			1100
1045			1101
1046	Xingtai Lv, Yuxin Zuo, Youbang Sun, Hongyi Liu, Yuntian Wei, Zhekai Chen, Lixuan He, Xuekai Zhu, Kaiyan Zhang, Bingning Wang, Ning Ding, and Bowen Zhou. 2025. Towards a unified view of large language model post-training . <i>CoRR</i> , abs/2509.04419.	Miguel Romero-Arjona, Pablo Valle, Juan C. Alonso, Ana Belén Sánchez, Miriam Ugarte, Antonia Cazalilla, Vicente Cambrón, José Antonio Parejo, Aitor Arrieta, and Sergio Segura. 2025. Red teaming contemporary AI models: Insights from spanish and basque perspectives . <i>CoRR</i> , abs/2503.10192.	1102
1047			1103
1048			1104
1049			1105
1050			1106
1051			1107
1052	Lu Ma, Hao Liang, Meiyi Qiang, Lexiang Tang, Xiaochen Ma, Zhen Hao Wong, Junbo Niu, Chengyu Shen, Runming He, Bin Cui, and Wentao Zhang. 2025. Learning what reinforcement learning can’t: Interleaved online fine-tuning for hardest questions . <i>CoRR</i> , abs/2506.07527.	Han Shen. 2025. On entropy control in LLM-RL algorithms . <i>CoRR</i> , abs/2509.03493.	1108
1053			1109
1054			1110
1055			1111
1056			1112
1057			1113
1058	Yingzhi Mao, Chunkang Zhang, Junxiang Wang, Xinyan Guan, Boxi Cao, Yaojie Lu, Hongyu Lin, Xianpei Han, and Le Sun. 2025. When models out-think their safety: Mitigating self-jailbreak in large reasoning models with chain-of-guardrails . <i>CoRR</i> , abs/2510.21285.	Parshin Shojaee, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. 2025. The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity . <i>Preprint</i> , arXiv:2506.06941.	1114
1059			1115
1060			1116
1061			1117
1062			1118
1063			1119
1064	Sara Vera Marjanović, Arkil Patel, Vaibhav Adlakha, Milad Aghajohari, Parishad BehnamGhader, Mehar	Viacheslav Sinii, Nikita Balagansky, Gleb Gerasimov, Daniil Laptev, Yaroslav Aksenov, Vadim Kurochkin, Alexey Gorbatovski, Boris Shaposhnikov, and Daniil Gavrilov. 2025a. Small vectors, big effects: A mechanistic study of rl-induced reasoning via steering vectors . <i>CoRR</i> , arXiv:2509.06608.	1120
1065			1121

1120	Viacheslav Sinii, Alexey Gorbatovski, Artem Cherepanov, Boris Shaposhnikov, Nikita Balagan-sky, and Daniil Gavrilov. 2025b. Steering llm reasoning through bias-only adaptation . <i>Preprint</i> , arXiv:2505.18706.	1174
1121		1175
1122		1176
1123		1177
1124		1178
1125	Linxin Song, Taiwei Shi, and Jieyu Zhao. 2025a. The hallucination tax of reinforcement finetuning . <i>CoRR</i> , abs/2505.13988.	1179
1126		1180
1127		1181
1128	Yuda Song, Julia Kempe, and Rémi Munos. 2025b. Outcome-based exploration for LLM reasoning . <i>CoRR</i> , abs/2509.06941.	1182
1129		1183
1130		1184
1131	Jinyan Su, Jennifer Healey, Preslav Nakov, and Claire Cardie. 2025a. Between underthinking and overthinking: An empirical study of reasoning length and correctness in llms . <i>CoRR</i> , abs/2505.00127.	1185
1132		1186
1133		1187
1134		1188
1135	Mingyu Su, Jian Guan, Yuxian Gu, Minlie Huang, and Hongning Wang. 2025b. Trust-region adaptive policy optimization . <i>Preprint</i> , arXiv:2512.17636.	1189
1136		1190
1137		1191
1138	Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Shaochen Zhong, Na Zou, Hanjie Chen, and Xia Hu. 2025. Stop overthinking: A survey on efficient reasoning for large language models . <i>Trans. Mach. Learn. Res.</i> , 2025.	1192
1139		1193
1140		1194
1141		1195
1142		1196
1143		1197
1144	Renliang Sun, Wei Cheng, Dawei Li, Haifeng Chen, and Wei Wang. 2025a. Stop when enough: Adaptive early-stopping for chain-of-thought reasoning . <i>CoRR</i> , abs/2510.10103.	1198
1145		1199
1146		1200
1147		1201
1148	Zhongxiang Sun, Qipeng Wang, Haoyu Wang, Xiao Zhang, and Jun Xu. 2025b. Detection and mitigation of hallucination in large reasoning models: A mechanistic perspective . <i>CoRR</i> , abs/2505.12886.	1202
1149		1203
1150		1204
1151		1205
1152	Yuqiao Tan, Minzheng Wang, Shizhu He, Huanxuan Liao, Chengfeng Zhao, Qiunan Lu, Tian Liang, Jun Zhao, and Kang Liu. 2025. Bottom-up policy optimization: Your language model policy secretly contains internal policies . <i>arXiv preprint arXiv:2512.19673</i> .	1206
1153		1207
1154		1208
1155		1209
1156		1210
1157		1211
1158	Miles Turpin, Andy Ardit, Marvin Li, Joe Benton, and Julian Michael. 2025. Teaching models to verbalize reward hacking in chain-of-thought reasoning . <i>CoRR</i> , abs/2506.22777.	1212
1159		1213
1160		1214
1161		1215
1162	Karthik Valmeekam, Kaya Stechly, Vardhan Palod, Atharva Gundawar, and Subbarao Kambhampati. 2025. Beyond semantics: The unreasonable effectiveness of reasonless intermediate tokens .	1216
1163		1217
1164		1218
1165		1219
1166	Constantin Venhoff, Iván Arcuschin, Philip Torr, Arthur Conmy, and Neel Nanda. 2025a. Base models know how to reason, thinking models learn when . <i>CoRR</i> , arXiv:2510.07364.	1220
1167		1221
1168		1222
1169		1223
1170	Constantin Venhoff, Iván Arcuschin, Philip Torr, Arthur Conmy, and Neel Nanda. 2025b. Understanding reasoning in thinking language models via steering vectors . <i>CoRR</i> , arXiv:2506.18167.	1224
1171		1225
1172		1226
1173		1227
		1228
	Chen Wang, Zhaochun Li, Jionghao Bai, Yuzhi Zhang, Shisheng Cui, Zhou Zhao, and Yue Wang. 2025a. Arbitrary entropy policy optimization: Entropy is controllable in reinforcement fine-tuning . <i>CoRR</i> , abs/2510.08141.	
	Haozhe Wang, Qixin Xu, Che Liu, Junhong Wu, Fangzhen Lin, and Wenhui Chen. 2025b. Emergent hierarchical reasoning in llms through reinforcement learning . <i>CoRR</i> , abs/2509.03646.	
	Jiayu Wang, Yifei Ming, Zixuan Ke, Caiming Xiong, Shafiq Joty, Aws Albarghouthi, and Frederic Sala. 2025c. Beyond accuracy: Dissecting mathematical reasoning for llms under reinforcement learning . <i>CoRR</i> , arXiv:2506.04723.	
	Kai Wang, Yihao Zhang, and Meng Sun. 2025d. When thinking llms lie: Unveiling the strategic deception in representations of reasoning models . <i>CoRR</i> , abs/2506.04909.	
	Yue Wang, Qiuzhi Liu, Jiahao Xu, Tian Liang, Xingyu Chen, Zhiwei He, Linfeng Song, Dian Yu, Juntao Li, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2025e. Thoughts are all over the place: On the underthinking of o1-like llms . <i>CoRR</i> , abs/2501.18585.	
	Zengzhi Wang, Fan Zhou, Xuefeng Li, and Pengfei Liu. 2025f. Octothinker: Mid-training incentivizes reinforcement learning scaling . <i>Preprint</i> , arXiv:2506.20512.	
	Jake Ward, Chuqiao Lin, Constantin Venhoff, and Neel Nanda. 2025a. Reasoning-finetuning repurposes latent representations in base models . <i>Preprint</i> , arXiv:2507.12638.	
	Jake Ward, Paul M. Riechers, and Adam Shai. 2025b. Rank-1 reasoning: Minimal parameter diffs encode interpretable reasoning signals . In <i>Mechanistic Interpretability Workshop at NeurIPS 2025</i> .	
	Fang Wu, Weihao Xuan, Ximing Lu, Zaïd Harchaoui, and Yejin Choi. 2025. The invisible leash: Why RLVR may not escape its origin . <i>CoRR</i> , abs/2507.14843.	
	Zhen Xiong, Yujun Cai, Zhecheng Li, and Yiwei Wang. 2025. Mapping the minds of llms: A graph-based analysis of reasoning LLM . <i>CoRR</i> , arXiv:2505.13890.	
	Fengli Xu, Qianye Hao, Zefang Zong, Jingwei Wang, Yunke Zhang, Jingyi Wang, Xiaochong Lan, Jiahui Gong, Tianjian Ouyang, Fanjin Meng, Chenyang Shao, Yuwei Yan, Qinglong Yang, Yiwen Song, Sijian Ren, Xinyuan Hu, Yu Li, Jie Feng, Chen Gao, and Yong Li. 2025. Towards large reasoning models: A survey of reinforced reasoning with large language models . <i>CoRR</i> , abs/2501.09686.	
	Ge Yan, Chung-En Sun, Tsui-Wei, and Weng. 2025. Reflectrl: Controlling llm reflection via representation engineering . <i>Preprint</i> , arXiv:2512.13979.	

1229	Shu Yang, Junchao Wu, Xin Chen, Yunze Xiao, Xinyi Yang, Derek F. Wong, and Di Wang. 2025a. Understanding aha moments: from external observations to internal mechanisms . <i>CoRR</i> , arXiv:2504.02956.	1285
1230		1286
1231		1287
1232		1288
1233	Wenkai Yang, Shuming Ma, Yankai Lin, and Furu Wei. 2025b. Towards thinking-optimal scaling of test-time compute for LLM reasoning . <i>CoRR</i> , abs/2502.18080.	1289
1234		1290
1235		1291
1236		1292
1237	Xinhao Yao, Lu Yu, Xiaolin Hu, Fengwei Teng, Qing Cui, Jun Zhou, and Yong Liu. 2025a. The debate on rlvr reasoning capability boundary: Shrinkage, expansion, or both? a two-stage dynamic view . <i>Preprint</i> , arXiv:2510.04028.	1293
1238		1294
1239		1295
1240		1296
1241		
1242	Yang Yao, Xuan Tong, Ruofan Wang, Yixu Wang, Lu-jundong Li, Liang Liu, Yan Teng, and Yingchun Wang. 2025b. A mousetrap: Fooling large reasoning models for jailbreak with chain of iterative chaos . In <i>Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025</i> , pages 7837–7855. Association for Computational Linguistics.	1297
1243		1298
1244		1299
1245		1300
1246		1301
1247		1302
1248		
1249		
1250	Zijun Yao, Yantao Liu, Yanxu Chen, Jianhui Chen, Jun-feng Fang, Lei Hou, Juanzi Li, and Tat-Seng Chua. 2025c. Are reasoning models more prone to hallucination? <i>CoRR</i> , abs/2505.23646.	1303
1251		1304
1252		1305
1253		1306
1254	Qingyu Yin, Chak Tou Leong, Linyi Yang, Wenxuan Huang, Wenjie Li, Xiting Wang, Jaehong Yoon, YunXing, XingYu, and Jinjin Gu. 2025. Refusal falls off a cliff: How safety alignment fails in reasoning? <i>CoRR</i> , abs/2510.06036.	1307
1255		1308
1256		1309
1257		
1258		
1259	Zonghao Ying, Guangyi Zheng, Yongxin Huang, Deyue Zhang, Wenxin Zhang, Quanchen Zou, Aishan Liu, Xianglong Liu, and Dacheng Tao. 2025. Towards understanding the safety boundaries of deepseek models: Evaluation and findings . <i>CoRR</i> , abs/2503.15092.	1310
1260		1311
1261		1312
1262		1313
1263		1314
1264	Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, and 16 others. 2025. Dapo: An open-source llm reinforcement learning system at scale . <i>Preprint</i> , arXiv:2503.14476.	1315
1265		1316
1266		1317
1267		1318
1268		
1269		
1270		
1271		
1272	Lifan Yuan, Weize Chen, Yuchen Zhang, Ganqu Cui, Hanbin Wang, Ziming You, Ning Ding, Zhiyuan Liu, Maosong Sun, and Hao Peng. 2025. From $f(x)$ and $g(x)$ to $f(g(x))$: LLMs learn new skills in rl by composing old ones . <i>Preprint</i> , arXiv:2509.25123.	1319
1273		1320
1274		1321
1275		1322
1276		1323
1277	Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. 2025. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? <i>Neural Information Processing Systems 39th (NeurIPS 2025)</i> , Oral Presentation.	1324
1278		1325
1279		
1280		
1281		
1282		
1283	Yuchen Zeng, Shuibai Zhang, Wonjun Kang, Shutong Wu, Lynnix Zou, Ying Fan, Heeju Kim, Ziqian Lin, Jungtaek Kim, Hyung Il Koo, Dimitris Papailiopoulos, and Kangwook Lee. 2025. Rejump: A tree-jump representation for analyzing and improving llm reasoning . <i>Preprint</i> , arXiv:2512.00831.	1326
1284		1327
		1328
		1329
		1330
		1331
	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 . <i>CoRR</i> , arXiv:2504.05419.	1332
		1333
		1334
		1335
		1336
		1337
	Charlie Zhang, Graham Neubig, and Xiang Yue. 2025b. On the interplay of pre-training, mid-training, and rl on reasoning language models . <i>Preprint</i> , arXiv:2512.07783.	
	Chong Zhang, Yue Deng, Xiang Lin, Bin Wang, Dianwen Ng, Hai Ye, Xingxuan Li, Yao Xiao, Zhanfeng Mo, Qi Zhang, and Lidong Bing. 2025c. 100 days after deepseek-r1: A survey on replication studies and more directions for reasoning language models . <i>CoRR</i> , abs/2505.00551.	
	Honglin Zhang, Qianyu Hao, Fengli Xu, and Yong Li. 2025d. Reinforcement learning fine-tuning enhances activation intensity and diversity in the internal circuitry of llms . <i>CoRR</i> , abs/2509.21044.	
	Jiajie Zhang, Nianyi Lin, Lei Hou, Ling Feng, and Juanzi Li. 2025e. Adaptthink: Reasoning models can learn when to think . <i>CoRR</i> , abs/2505.13417.	
	Jue Zhang, Qingwei Lin, Saravan Rajmohan, and Dongmei Zhang. 2025f. From reasoning to answer: Empirical, attention-based and mechanistic insights into distilled deepseek r1 models . <i>Preprint</i> , arXiv:2509.23676.	
	Junyu Zhang, Yifan Sun, Tianang Leng, Jingyan Shen, Liu Ziyin, Paul Pu Liang, and Huan Zhang. 2025g. When reasoning meets its laws . <i>arXiv preprint arXiv:2512.17901</i> .	
	Kaiyan Zhang, Yuxin Zuo, Bingxiang He, Youbang Sun, Runze Liu, Che Jiang, Yuchen Fan, Kai Tian, Guoli Jia, Pengfei Li, Yu Fu, Xingtai Lv, Yuchen Zhang, Sihang Zeng, Shang Qu, Haozhan Li, Shijie Wang, Yuru Wang, Xinwei Long, and 20 others. 2025h. A survey of reinforcement learning for large reasoning models . <i>CoRR</i> , abs/2509.08827.	
	Wenhao Zhang, Yuexiang Xie, Yuchang Sun, Yanxi Chen, Guoyin Wang, Yaliang Li, Bolin Ding, and Jingren Zhou. 2025i. On-policy RL meets off-policy experts: Harmonizing supervised fine-tuning and reinforcement learning via dynamic weighting . <i>CoRR</i> , abs/2508.11408.	
	Wenjing Zhang, Xuejiao Lei, Zhaoxiang Liu, Limin Han, Jiaojiao Zhao, Beibei Huang, Zhenhong Long, Junting Guo, Meijuan An, Rongjia Du, Ning Wang, Kai Wang, and Shiguo Lian. 2025j. Safety evaluation and enhancement of deepseek models in chinese contexts . <i>CoRR</i> , abs/2503.16529.	

1338	Haoran Zhao, Yuchen Yan, Yongliang Shen, Haolei Xu,
1339	Wenqi Zhang, Kaitao Song, Jian Shao, Wei-ming
1340	Lu, Jun Xiao, and Yueting Zhuang. 2025a. Let llms
1341	break free from overthinking via self-braking tuning.
1342	<i>CoRR</i> , abs/2505.14604.
1343	James Xu Zhao, Bryan Hooi, and See-Kiong Ng. 2025b.
1344	Test-time scaling in reasoning models is not ef-
1345	fective for knowledge-intensive tasks yet. <i>CoRR</i> ,
1346	abs/2509.06861.
1347	Weixiang Zhao, Xingyu Sui, Jiahe Guo, Yulin Hu, Yang
1348	Deng, Yanyan Zhao, Bing Qin, Wanxiang Che, Tat-
1349	Seng Chua, and Ting Liu. 2025c. Trade-offs in large
1350	reasoning models: An empirical analysis of delibera-
1351	tive and adaptive reasoning over foundational capa-
1352	bilities. <i>CoRR</i> , abs/2503.17979.
1353	Kaiwen Zhou, Chengzhi Liu, Xuandong Zhao, Shreed-
1354	har Jangam, Jayanth Srinivasa, Gaowen Liu, Dawn
1355	Song, and Xin Eric Wang. 2025. The hidden risks of
1356	large reasoning models: A safety assessment of R1.
1357	<i>CoRR</i> , abs/2502.12659.
1358	Rongzhi Zhu, Yi Liu, Zequn Sun, Yiwei Wang, and
1359	Wei Hu. 2025. When can large reasoning models
1360	save thinking? mechanistic analysis of behavioral
1361	divergence in reasoning. <i>CoRR</i> , abs/2505.15276.

A Taxonomy 1362

We present the taxonomy of our paper in Figure 2. We follow Figure 1 to organize the research of various directions and list representative works accordingly. 1363 1364 1365 1366

B Future Directions 1367

B.1 Applied Interpretability 1368

Mechanistic interpretability (MI) research is increasingly illuminating the internal logic of LRMs. A crucial next step is to leverage these insights for targeted improvements, moving from passive understanding to active application. 1369 1370 1371 1372 1373

Training-Time Applications. A promising direction lies in using internal representations to directly inform RL algorithm design. Recent studies demonstrate initial success in this area, such as utilizing attention mechanisms to inform reward shaping (Li et al., 2025e) or policy sampling strategies (Liu et al., 2025c), and decoding intermediate layer activations to infer latent policies (Tan et al., 2025). The overarching challenge is to systematically transform mechanistic insights into algorithmic improvements for RL components. 1374 1375 1376 1377 1378 1379 1380 1381 1382 1383 1384

Inference-Time Applications. Mechanistic findings can also be applied to steer model behavior during inference. While existing work already uses insights of reasoning structures or specific representations to improve performance, deeper opportunities remain. For instance, research suggests that RL may not effectively leverage the full depth of models (Hu et al., 2025b; Nepal et al., 2025). This understanding should actively inform the design of novel training algorithms and architectures that better utilize internal computational pathways. 1385 1386 1387 1388 1389 1390 1391 1392 1393 1394 1395

B.2 Advancing Interpretability Methodology 1396

Future research should emphasize developing scalable and generalizable MI frameworks specifically tailored for LRMs. First, the enormous scale of LRMs in *training cost*, *inference length*, and *parameter count* poses significant methodological challenges. Conducting controlled experiments to isolate variables is difficult, and techniques like training SAEs become computationally prohibitive, slowing progress and reducing reproducibility. There is a clear need for more efficient and scalable MI tools tailored to these models. Second, many MI findings remain model-specific, failing to generalize across different architectures or 1397 1398 1399 1400 1401 1402 1403 1404 1405 1406 1407 1408 1409

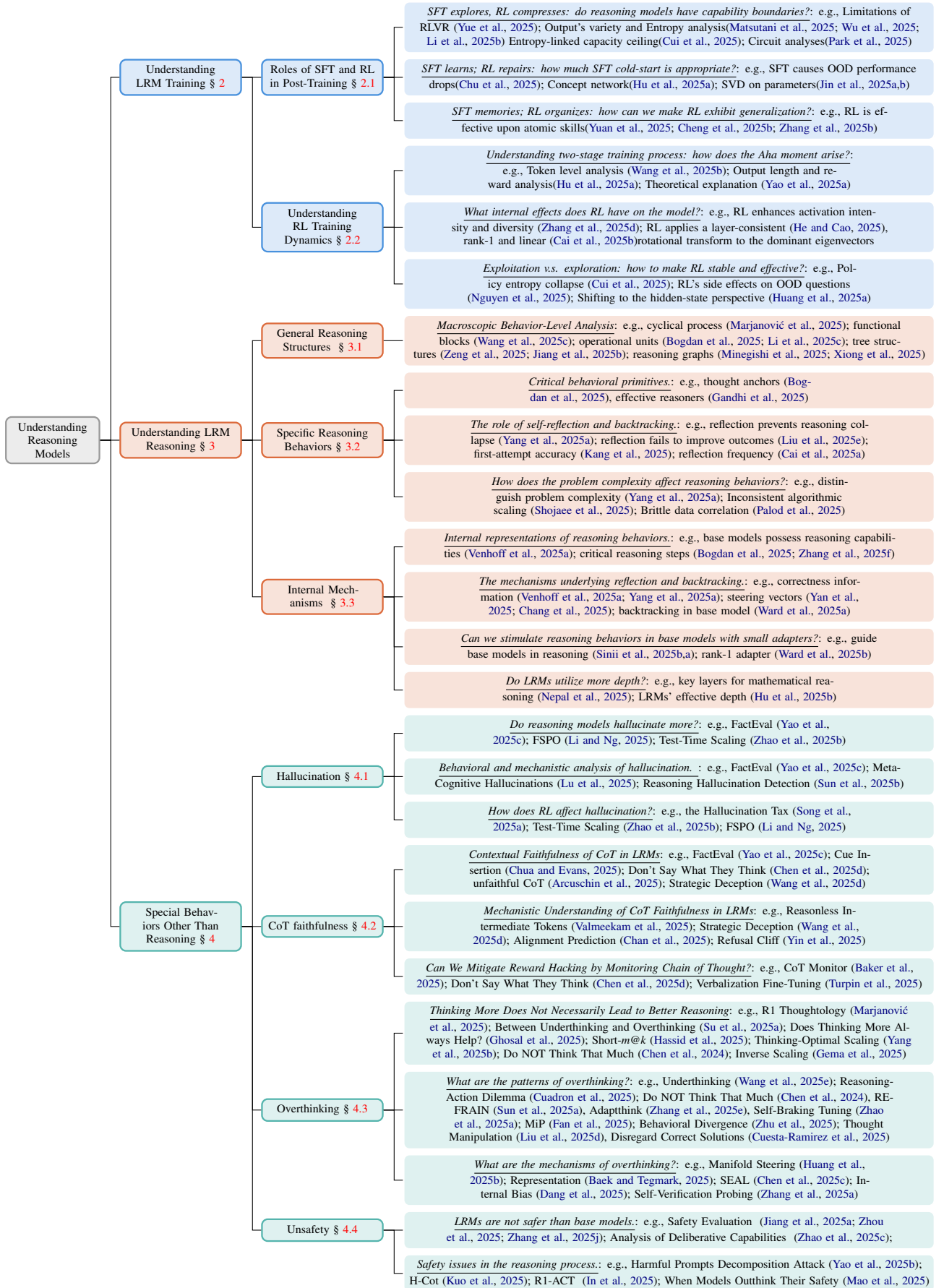


Figure 2: Taxonomy of our paper and representative works for each direction.

1410 training runs. To enhance scientific value, the field
1411 should strive for general frameworks that abstract
1412 away implementation details and uncover universal
1413 reasoning principles. This could involve establish-
1414 ing benchmarks for mechanistic generalization, de-
1415 veloping theory-grounded methods less sensitive to
1416 model quirks, or building more robust interpretabil-
1417 ity probes.

1418 **B.3 Toward A Unified Theory**

1419 Current mechanistic research has produced a
1420 wealth of empirical findings—model-specific pat-
1421 terns, dataset-specific behaviors, and localized ex-
1422 planations for special phenomena. Yet it lacks a
1423 predictive, fundamental science of reasoning in
1424 LRMs. Such a theory should be *fundamental*, ab-
1425 stracting from implementation to reveal first prin-
1426 ciples governing reasoning. Early efforts, such as
1427 theoretically formalizing laws of reasoning (Zhang
1428 et al., 2025g) that link task complexity to model be-
1429 havior, mark a step in this direction. Furthermore, a
1430 mature theory should be *predictive*. Similar to scal-
1431 ing laws in LLM pre-training (Kaplan et al., 2020),
1432 it should forecast model behaviors and to establish
1433 a set of “laws of reasoning” that not only explain
1434 existing empirical results but also actively guide
1435 the design of future models, training algorithms,
1436 and evaluation frameworks, transforming MI from
1437 a descriptive tool into a foundational science.

1438 **C Training Methods**

1439 **C.1 Combine SFT with RL**

1440 Running SFT and RL as two separate steps will
1441 let the bias introduced by SFT grow too large and
1442 degrade final performance. Therefore, some stud-
1443 ies attempt to combine the two approaches into a
1444 unified single post-training step.

1445 Some explorations primarily focused on inter-
1446 leaving the SFT and RL processes and on identi-
1447 fying appropriate switching points between them.
1448 Based on their research into RL training dynamics,
1449 Hu et al. (2025a) proposed the Annealed-RLVR
1450 algorithm, which introduces SFT for heating when
1451 accuracy is very low to disrupt the current subopti-
1452 mal state, then continues RL to perform annealing.
1453 Ma et al. (2025) observes that RL excels at easy
1454 questions while SFT is better suited to hard ones;
1455 their *ReLIFT* pipeline automatically flags the hard
1456 instances during RL, collects corresponding expert
1457 demonstrations, and inserts an SFT update once
1458 enough difficult question–answer pairs have been

1459 accumulated. *TRAPO* (Su et al., 2025b) interleaves
1460 SFT and RL within every training instance and sets
1461 up a mechanism that dynamically supplies expert-
1462 guided prefixes. SFT in *TRAPO* is constrained by
1463 trust-region gradient clipping to avoid distribution-
1464 blending.

1465 Further more, some studies aim to fuse the loss
1466 functions of RL and SFT to achieve a truly unified
1467 post-training approach. *UFT* (Liu et al., 2025b)
1468 introduces an additional log-likelihood term to the
1469 objective function of RFT(RL Fine-Tuning), al-
1470 lowing the model to learn from the informative
1471 supervision signal and still benefit from the gen-
1472 eralization of RFT. *HPT* (Lv et al., 2025) defines
1473 the total loss as a weighted sum of the SFT and RL
1474 losses and dynamically adjusts the weights of SFT
1475 and RL based on real-time performance. *CHORD*
1476 (Zhang et al., 2025i) treats SFT as a dynamically-
1477 weighted auxiliary objective within the RL process
1478 and introduces a token-level weighting function
1479 that up-weights the SFT component only when the
1480 model is uncertain about the answer. Going further,
1481 *SRFT* (Fu et al., 2025) incorporates demonstration
1482 data into the RL training set and constructs the fi-
1483 nal loss as the entropy-weighted sum of four terms:
1484 SFT loss on demonstrations, RL loss on demon-
1485 strations, and RL loss on positive (entropy-weighted)
1486 and negative (without weighting) sampled rollouts.
1487 *BRIDGE* (Chen et al., 2025a) attaches LoRA fine-
1488 tuning blocks to the model architecture and for-
1489 mulates the post-training procedure as a bi-level
1490 optimization: SFT refines the model starting from
1491 the parameter optimum found by RL, optimizing
1492 solely over the LoRA weights, with an appropri-
1493 ate transformation eliminates the need for second-
1494 order gradients.

1495 **C.2 RL balancing exploration and 1496 exploitation**

1497 To address entropy collapse and balance explo-
1498 ration and exploitation for stable RL training, nu-
1499 merous studies have proposed solutions. *DAPO*
1500 (Yu et al., 2025) decouples the clipping bounds
1501 of PPO into ϵ -high and ϵ -low, raising ϵ -high to
1502 leave more head-room for boosting the probabili-
1503 ties of low-probability “exploratory” tokens. (Cui
1504 et al., 2025) shows that the entropy change is gov-
1505 erned by the covariance between the “action log-
1506 probabilities” and the “changes in action logits”;
1507 tokens with high covariance are the main drivers
1508 of entropy collapse. To counter this, they propose
1509 *Clip-Cov* which randomly truncates the gradients of

1510 high-covariance tokens, and *KL-Cov* which adds an
1511 extra KL-penalty to those tokens. *ProRL* (Liu et al.,
1512 2025a) adopts the Decoupled Clip technique from
1513 *DAPO*, and further equips the pipeline with KL
1514 regularization plus periodic reference-policy resets
1515 to avert entropy collapse, enabling effective RL
1516 training that continues for thousands of steps.

1517 More recent studies have moved beyond simple
1518 clipping and experimented with additional mecha-
1519 nisms to further arrest entropy collapse. *CURE* (Li
1520 et al., 2025d) shows that the prevailing RLVR
1521 pipeline relies on sampling from a fixed initial state,
1522 biasing the model toward overly deterministic be-
1523 havior and low diversity. They introduce a two-
1524 stage scheme to balance exploration and exploita-
1525 tion: in stage one they identify high-entropy tokens,
1526 truncate at those tokens, and then sample multiple
1527 continuations that are all used for updating, thereby
1528 intensifying exploration around high-entropy re-
1529 gions; in stage two they revert to the ordinary static-
1530 sampling *DAPO* routine. Similarly, Nguyen et al.
1531 (2025) advocates sampling problems that the base
1532 model still handles poorly, rather than repeatedly
1533 drawing those it already solves well. Song et al.
1534 (2025b) encourages historical exploration of rare
1535 answers through UCB-style rewards and fosters
1536 batch exploration at test time by penalizing dupli-
1537 cate answers within each batch.

1538 Naive entropy regularization performs poorly
1539 when training reasoning models. Several works
1540 have designed regularization methods specifically
1541 tailored to reasoning models. Cheng et al. (2025a)
1542 injects a clipped, gradient-detached entropy term
1543 into the advantage function to encourage longer
1544 chains-of-thought. Wang et al. (2025a) stabilizes
1545 policy entropy by combining Policy-Gradient, Dis-
1546 tribution, and Reinforce signals into a composite
1547 regularizer. Shen (2025); Jiang et al. (2025c) com-
1548 pute entropy only over the top-p tokens and adap-
1549 tively rescale it for entropy regularization; the lat-
1550 ter research further shows that regularizing those
1551 high-entropy tokens only can improve model per-
1552 formance.