

# Opening the Black Box: A Survey on the Mechanisms of Multi-Step Reasoning in Large Language Models

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) have demonstrated remarkable abilities to solve problems requiring multiple reasoning steps, yet the internal mechanisms enabling such capabilities remain elusive. Unlike existing surveys that primarily focus on engineering methods to enhance performance, this survey provides a comprehensive overview of the *mechanisms* underlying LLM multi-step reasoning. We organize the survey around a conceptual framework comprising seven interconnected research questions from how LLMs execute *implicit multi-hop reasoning* within hidden activations to how *verbalized explicit reasoning* remodels the internal computation. Finally, we highlight five research directions for future mechanistic studies.

## 1 Introduction

Large Language Models (LLMs) have demonstrated an impressive ability to carry out *multi-step reasoning*, which involves the process of drawing conclusions through a sequence of intermediate steps, where each step builds on the previous one. Multi-step reasoning has been widely regarded as one of the most fundamental forms of reasoning (Hou et al., 2023a; Guo et al., 2025). It serves as the backbone of advanced tasks such as deep question answering, mathematical problem solving, logical deduction, code generation, and planning (OpenAI, 2023; Yang et al., 2024a; Guo et al., 2024; DeepSeek-AI et al., 2025).

Multi-step reasoning in LLMs generally takes on two distinct forms. *Implicit reasoning* involves performing multi-hop inference entirely within the model’s hidden activations, delivering a correct final answer without verbalizing intermediate steps. In contrast, *explicit reasoning*, exemplified by *Chain-of-Thought* (CoT) (Wei et al., 2022b), instructs the model to externalize the reasoning process into a sequence of natural language tokens. Remarkably, modern LLMs have exhibited strong

performance in both paradigms (Chu et al., 2024a; Chen et al., 2025a; Li et al., 2025b). Building on this empirical success, the *internal mechanisms* that enable such capabilities become scientifically intriguing. For implicit reasoning, a key puzzle is *how* multi-step reasoning capabilities emerge from simple next-token prediction training, and *how* LLMs internally carry out multi-step computations. For explicit CoT reasoning, critical questions persist about *why* CoT can elicit superior reasoning capabilities and *whether* the generated rationale faithfully reflects the model’s actual decision-making process. Understanding these mechanisms is not only a matter of scientific curiosity but also a prerequisite for building more reliable, controllable, and human-aligned reasoning systems.

Although we still lack a unified mechanistic theory, a growing body of literature seeks to *open the black box of LLM multi-step reasoning* and has made significant progress. In this paper, we aim to provide a comprehensive overview of these works. Unlike existing surveys (Huang and Chang, 2023; Chu et al., 2024b; Chen et al., 2025a) that primarily focus on *enhancing* reasoning (*e.g.*, through tool use, retrieval augmentation, or self-correction), our survey explicitly focuses on *understanding mechanisms*, a perspective that has been largely overlooked in previous reviews. As illustrated in Figure 1, we identify seven pivotal, interconnected, and progressive *research questions (RQs)* to form the cognitive framework of our survey. These questions form a cohesive narrative, covering analytical methods and key findings from the hidden internal dynamics of latent reasoning to the visible mechanisms of explicit CoT reasoning. We end by pointing out five promising future research directions.

## 2 Implicit Multi-Step Reasoning

Implicit multi-step reasoning involves answering questions by synthesizing information across mul-

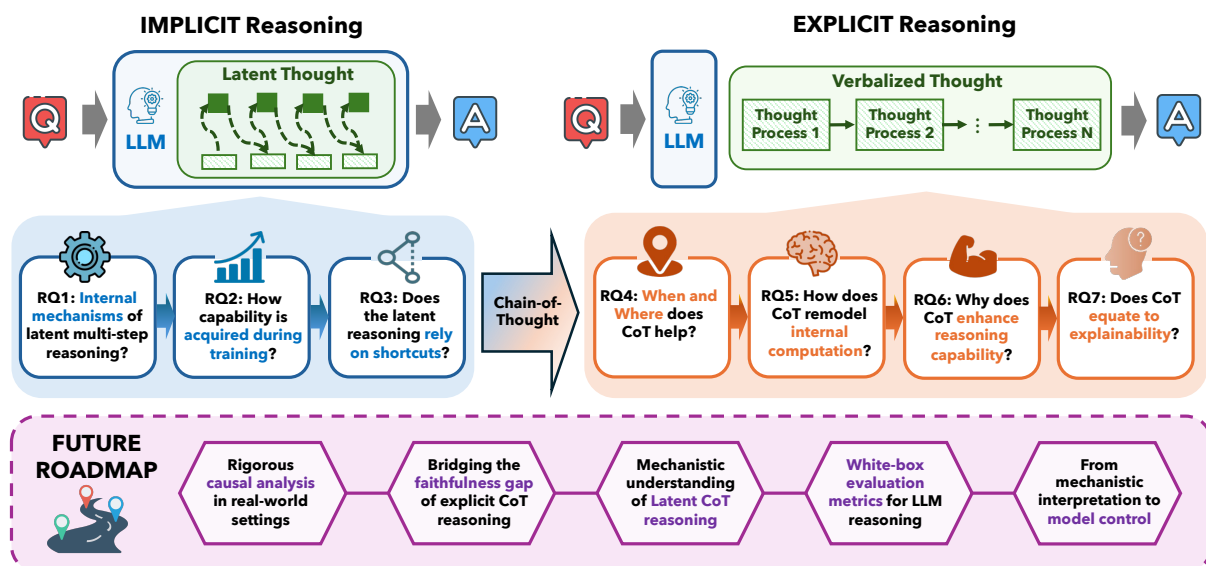


Figure 1: The cognitive framework and organizational structure of this survey. We explore the mechanisms of multi-step reasoning through two distinct paradigms: *Implicit Reasoning* and *Explicit Reasoning*, through seven interconnected *Research Questions*. The bottom panel highlights five strategic directions for future research.

081 tiple steps entirely within hidden states, without  
 082 verbalizing intermediate steps. Understanding the  
 083 mechanisms behind this capability is critical to  
 084 determine whether models perform genuine step-  
 085 by-step inference or rely on shallow shortcuts, ulti-  
 086 mately guiding the development of more trustwor-  
 087 thy and generalizable reasoning models.

## 088 2.1 What are the internal mechanisms of 089 latent multi-step reasoning?

090 Recent mechanistic studies have begun to unveil  
 091 how LLMs carry out latent multi-hop computation  
 092 entirely in their hidden states (Yang et al., 2024b;  
 093 Biran et al., 2024; Brinkmann et al., 2024a). These  
 094 studies employing causal probing, mechanistic trac-  
 095 ing, and representational analysis have collectively  
 096 revealed a *staged* internal process in which interme-  
 097 diate results are computed and transformed layer  
 098 by layer, ultimately contributing to the final output.  
 099 In essence, transformers appear to implement an  
 100 internal chain-of-thought spread across their depth.

101 **Functional specialization of layers.** A major  
 102 body of work explores *layer specialization*, aiming  
 103 to identify the distinct computational roles each  
 104 layer plays during multi-hop inference. Using  
 105 *Patchscopes* (Ghandeharioun et al., 2024) together  
 106 with a novel intervention technique termed *back-*  
 107 *patching*, Biran et al. (2024) uncovered a sequen-  
 108 tial computation pathway in which early layers iden-  
 109 tify the bridge entity, which is then propagated  
 110 forward and exploited by later layers to complete

the inference. Complementarily, Li et al. (2024a)  
 applied logit lens analysis (nostalgebraist, 2020)  
 and found that implicit reasoning representations  
 emerge in intermediate layers and have a causal  
 influence on generating the final answer. Extend-  
 ing this perspective, Yu et al. (2025) traced log-  
 its through the network via a neuron-level logit  
 flow method and observed that even a single-hop  
 query is solved in multiple distinct stages—entity  
 subject enrichment, entity attribute extraction, re-  
 lation subject enrichment, and relation attribute  
 extraction—each of which is localized to different  
 layers. More recently, Yang et al. (2025d) showed  
 that this layer-wise reasoning also applies at the  
 task level: for composite instructions, models exe-  
 cute different subtasks at different depths, forming  
 a staged computation across layers. All the above  
 studies provided evidence of functional specializa-  
 tion of transformer layers in multi-hop reasoning.

130 **Uncovering fine-grained reasoning structures.**  
 131 Beyond layer specification, another line of work  
 132 aims to recover more fine-grained implicit rea-  
 133 soning structures from model internals. *Mech-*  
 134 *anisticProbe* (Hou et al., 2023b) introduced an  
 135 attention-probing technique to extract latent rea-  
 136 soning trees from transformer activations. They  
 137 showed that on synthetic and natural tasks with  
 138 GPT-2 and LLaMA, models often perform proce-  
 139 dural reasoning layer by layer, with lower layers  
 140 selecting statements and higher layers executing  
 141 reasoning steps. Complementing these findings,

142 Brinkmann et al. (2024b) analyzed a small trans- 193  
143 former trained on a symbolic tree path-finding task, 194  
144 finding that it implements a backward chaining algo- 195  
145 rithm: deduction heads climb trees one level per 196  
146 layer, register tokens act as working memory for 197  
147 parallel subpaths, and a one-step lookahead heuris- 198  
148 tic compensates when chaining is insufficient. To- 199  
149 gether, these studies demonstrate that transform- 200  
150 ers can adopt structured, algorithm-like reasoning 201  
151 strategies beyond memorization, albeit within the 202  
152 limits of the model’s depth (to be discussed below). 203

153 **Layer depth as the primary bottleneck for im- 204**  
154 plicit reasoning. Theoretical and empirical stud- 205  
155 ies indicate that the number of reasoning steps a 206  
156 model can perform implicitly is strictly limited by 207  
157 its depth. Merrill and Sabharwal (2024) theoret- 208  
158 ically demonstrated that a standard Transformer 209  
159 with constant depth cannot solve inherently serial 210  
160 problems that require computation scaling with 211  
161 input size, *e.g.*, parity or graph connectivity. In 212  
162 practice, Yu (2025) and Guo et al. (2025) found 213  
163 that specific multi-hop reasoning tasks require a 214  
164 minimum threshold of layers to resolve; if a model 215  
165 is too shallow, the “latent chain” is cut short, and 216  
166 the reasoning fails. Saunshi et al. (2025) formally 217  
167 established that an  $L$ -layer Transformer can simu- 218  
168 late an  $m$ -step explicit reasoning process, provided 219  
169  $L$  is sufficiently large to accommodate the iterative 220  
170 forward passes required. All these works revealed 221  
171 a close correlation between layer depth and the 222  
172 implicit reasoning capabilities of the model. 223

173 **Why implicit reasoning sometimes fails.** Ident- 224  
174 ifying how and why implicit reasoning sometimes 225  
175 fails has also been illuminating. Biran et al. (2024) 226  
176 discovered that many failures stem from delayed 227  
177 resolution of the first hop, and showed that rerun- 228  
178 ning computations via back-patching can correct 229  
179 these errors. Li et al. (2024a) found that failures 230  
180 frequently arise from the improper generation or 231  
181 utilization of implicit reasoning results. To address 232  
182 this, they proposed CREME, a lightweight model- 233  
183 editing technique that patches specific multi-head 234  
184 self-attention modules, leading to improved com- 235  
185 positional reasoning generalization with minimal 236  
186 disruption to unrelated predictions. In the context 237  
187 of two-hop queries (“ $e_1$ ’s  $r_1$ ’s  $r_2$  is  $e_3$ ”), Yu et al. 238  
188 (2025) showed that errors often occur when high- 239  
189 layer features at the  $r_1$  position overemphasize the 240  
190 intermediate entity  $e_2$ , outweighing the logits for 241  
191 the correct final entity  $e_3$ . This finding revealed that 242  
192 LLMs internally build and combine entity–relation

representations in a staged manner, but positional 193  
interference can derail multi-hop reasoning. To 194  
fix this, they introduced a back-attention mecha- 195  
nism allowing lower layers to reuse higher-layer 196  
information from other positions, which substan- 197  
tially improved multi-hop accuracy. However, even 198  
with such interventions, certain transformers still 199  
struggle to reliably chain more than one reasoning 200  
step. For example, Yang et al. (2024b) found that 201  
*LLaMA-2* models, while reliably recalling a needed 202  
bridge entity, often fail to apply it to the second 203  
hop, highlighting limits in architecture that impede 204  
consistent multi-step chaining. 205

**Takeaway:** Implicit multi-hop reasoning in 206  
LLMs is not a monolithic capability but rather 207  
an orchestrated, layered process, with distinct 208  
modules and pathways specializing in differ- 209  
ent phases of the reasoning chain. For exam- 210  
ple, probing and intervention studies showed 211  
that intermediate results, *e.g.*, bridge entities, 212  
are computed and passed along inside the net- 213  
work. Nevertheless, such implicit reasoning 214  
is constrained by the inherent architecture of 215  
transformers, for example, their fixed depth. 216

## 2.2 How latent multi-step reasoning capability 207 is acquired during training? 208

209 Models do not possess latent reasoning capabilities 210  
at initialization. If multi-hop reasoning is imple- 211  
mented via specialized internal circuits discussed 212  
in Section 2.1, a critical question arises: *how do* 213  
*these circuits emerge in the first place?* Research 214  
into training dynamics reveals that implicit reason- 215  
ing is an *acquired* behavior that emerges during the 216  
training process through distinct phase transitions.

217 **Grokking marks the shift from memorization to 218**  
219 reasoning. Recent studies (Wang et al., 2024; Ye 220  
221 et al., 2025; Zhang et al., 2025c; Abramov et al., 222  
223 2025) suggested that LLMs do not learn multi-step 224  
225 reasoning gradually; instead, they often undergo 226  
227 *phase transitions* during training where reasoning 228  
229 capabilities appear *suddenly* rather than continu- 230  
231 ously. In other words, a model might spend many 232  
updates seemingly memorizing or floundering, then 233  
“grok” the underlying reasoning algorithm after 234  
a certain point. This phenomenon, known as 235  
“*grokking*”, was initially observed in deep networks 236  
trained on other tasks such as modular arithmetic, 237  
where generalization performance spikes long after 238  
training accuracy has saturated (Power et al., 2022; 239  
240 241

Olsson et al., 2022; Wei et al., 2022a).

In the context of multi-hop implicit reasoning, this phenomenon of transformers transitioning from early-stage *memorization* to later-stage *generalization* was first observed by Wang et al. (2024) through training transformers from scratch on symbolic reasoning tasks. They found that the multi-hop reasoning capability emerges only through *grokking*, where an early *memorizing circuit* is gradually replaced by a more efficient *generalizing circuit* due to optimization bias and weight decay. Ye et al. (2025) corroborated this phase transition, proposing a *three-stage trajectory*: (i) rapid memorization, (ii) delayed in-distribution generalization, and (iii) slower cross-distribution generalization, with persistent OOD bottlenecks at the second hop. Mechanistically, they employed *cross-query semantic patching* to localize the “bridge” entity and a *cosine-based representational lens* to reveal that generalization coincides with *mid-layer clustering* of intermediate entity representations.

### Factors influencing the emergence of reasoning.

The transition from memorization to generalization is not random; studies revealed that it is governed by specific properties. One of the primary determinants is the *training data distribution*. Wang et al. (2024) demonstrated that the speed of grokking correlates strongly with the *ratio of inferred to atomic facts*  $\phi$  in training. A higher ratio of compositional examples forces the model to abandon inefficient memorization in favor of the generalizing circuit. Expanding this to real-world scenarios, Abramov et al. (2025) found that natural corpora often lack sufficient connectivity (low  $\phi$ ) to trigger this transition, but data augmentation with synthetic inferred facts can artificially raise  $\phi$  above the critical threshold required for circuit formation. Beyond data distribution, the *scale of the training data* also matters. Yao et al. (2025b) revealed a scaling law: the data budget required to learn implicit  $k$ -hop reasoning grows exponentially with  $k$ , though curriculum learning can significantly mitigate this cost. From an optimization perspective, Zhang et al. (2025c) identified *complexity control* parameters as crucial factors. They found that smaller initialization scales and stronger weight decay bias the optimization process toward low-complexity, rule-like solutions rather than high-complexity, memory-based mappings, thereby accelerating the emergence of reasoning capabilities. Finally, Li et al. (2025c) observed that in large-

scale pretraining, grokking is *asynchronous and local*; different domains and data groups undergo this memorization-to-generalization transition at different times depending on their inherent difficulty and distribution heterogeneity.

**Takeaway:** Implicit multi-hop reasoning capability is an *acquired* capability that emerges via *grokking*—a phase transition from surface-level memorization to structured reasoning. This transition is not automatic; it is governed by critical factors, including the training data distribution, the data scale, and complexity control via optimization biases.

### 2.3 To what extent does multi-step reasoning rely on shortcuts?

While the training dynamics discussed in § 2.1 suggest that structured reasoning circuits can emerge, growing mechanistic evidence has also uncovered a more complex and often discouraging reality regarding model internals. Models frequently bypass genuine multi-step reasoning, relying instead on “*shortcuts*”—statistical correlations or surface-level heuristics that mimic reasoning without performing the underlying computation (Kang and Choi, 2023; Elazar et al., 2024; Yang et al., 2025b).

#### Factual shortcuts bypass intermediate reasoning.

A primary form of shortcutting involves exploiting direct associations between the subject and the final answer, effectively skipping the intermediate steps. Ju et al. (2024) investigated this in the context of knowledge editing, finding that failures often stem from “*shortcut neurons*” that encode a direct link between the first and last entities, ignoring the multi-hop structure. Mechanistically, Yang et al. (2025c) used *Patchscopes* (Ghandeharioun et al., 2024) to distinguish valid reasoning from shortcuts. They observed that genuine implicit reasoning coincides with the model constructing a hidden representation of the intermediate bridge entity. In contrast, shortcut-prone queries bypass this internal construction entirely. When these direct shortcuts are removed, model performance drops by nearly a factor of three, revealing that much of the perceived reasoning capability is illusory.

#### Shortcuts based on surface-level pattern matching.

Beyond factual associations, models also latch onto structural regularities in the training data. Lin et al. (2025) analyzed implicit arithmetic rea-

soning and found that models often adopt a “bag-of-words” heuristic, treating operations as commutative even when they are not. While this shortcut works for fixed-template examples, performance collapses when premise order is randomized, proving the model had not learned the robust sequential logic. Similarly, Guo et al. (2025) found that in the presence of context distractors, pretrained models default to a heuristic of guessing based on surface plausibility. However, they also noted a positive trajectory: fine-tuning can force a phase transition where the model shifts from this shallow guessing behavior to a sequential query mechanism that explicitly retrieves intermediate entities.

**Takeaway:** LLMs frequently bypass the “latent reasoning chain” via *factual shortcuts* (direct input-output associations) or *structural heuristics* (exploiting surface patterns like commutativity). This underscores the need for shortcut-free evaluation protocols and training setups that force models to construct and reuse intermediate representations.

### 3 Explicit Multi-Step Reasoning

Implicit reasoning operates entirely within the fixed computational budget of the model’s hidden states; therefore, it is bounded by the depth bottleneck and frequently falls prey to shortcuts. *Explicit multi-step reasoning* fundamentally alters this paradigm. By prompting an LLM to produce a step-by-step *Chain-of-Thought* (CoT), the reasoning process is externalized into a sequence of natural language tokens, effectively extending the computational capacity beyond the model’s layers. CoT has been shown to unlock significantly better performance on tasks that require reasoning. In this section, we dissect the mechanisms of this paradigm through four progressive research questions (§ 3.1-§ 3.4).

#### 3.1 Where and When Does CoT Help?

**On which tasks does CoT help?** To uncover this, Sprague et al. (2025) conducted a large-scale meta-analysis across 20 benchmarks and found that prompting with CoT yields large gains primarily on *math and symbolic logic tasks*, with far smaller or even negative gains on other domains. Suzgun et al. (2023) similarly showed that many *BIG-Bench Hard* tasks (Srivastava et al., 2023), which had stumped standard few-shot prompts, become solvable with CoT. These were precisely

tasks requiring multi-step reasoning, e.g., symbolic manipulation, compositional logic. However, for knowledge-heavy tasks like MMLU (Hendrycks et al., 2021) or commonsense reasoning, CoT often provides negligible improvement (Sprague et al., 2025). In certain cases, CoT can even degrade accuracy. For example, Liu et al. (2024) examined cognitive-psychology tasks where additional deliberation harms human performance, e.g., certain trick riddles or intuitive judgment problems. They found that CoT substantially degraded accuracy on such tasks, and it tends to distract the model into over-complicating a problem that might have been solved via intuition. A complementary study on Blocksworld planning (Stechly et al., 2024) found that CoT helps only when the prompt examples closely match the test distribution, and the gains quickly deteriorate if the test problem’s complexity exceeds that seen in the exemplars.

**What factors influence the efficacy of CoT?** Beyond task-level evaluations, empirical studies have shown that CoT performance can be dramatically influenced by many features of the CoT prompt. First, studies (Ye and Durrett, 2022; Madaan et al., 2023; Wang et al., 2023) reveal that the *relevance and ordering* of exemplars matter more than their semantics; models can still derive correct answers from invalid rationales if the prompt maintains a coherent structure. Second, the length of reasoning is another critical factor, with Jin et al. (2024) identifying that the number of reasoning steps significantly modulates model performance. Finally, CoT is surprisingly sensitive to phrasing; minor input perturbations can substantially bias models’ answers (Turpin et al., 2023; Sadr et al., 2025).

**Why do these factors influence CoT efficacy?** To explain the mechanisms underlying these factors, recent research provided theoretical and mechanistic groundings. Tutunov et al. (2023) proposed that CoT efficacy stems from the model’s ability to approximate the true conditional distribution of reasoning, where structured exemplars help the model infer the task’s latent logic and reduce generation ambiguity. Prabhakar et al. (2024) refined this view through a controlled case study, characterizing CoT as a probabilistic process heavily modulated by output *probability*, task *memorization* in training data, and step-wise *complexity*. Mechanistically, Wu et al. (2023) revealed how specific components of the CoT prompt drive model generation via gradient-based feature attribution.

**Takeaway:** CoT prompting yields significant gains primarily in tasks involving *mathematical, logical, or symbolic reasoning*. Its efficacy depends more on the structural coherence and relevance of exemplars, the length of reasoning, and the prompt phrasing. Several theoretical and mechanistic frameworks were proposed to understand such driving factors.

### 3.2 How Does Chain-of-Thought Remodel Internal Computation?

Chain-of-thought prompting does more than just alter an LLM’s output format. Growing evidence shows that it fundamentally changes the model’s internal computation into a “reasoning mode”, where the model retrieves and updates information in a stepwise fashion, leveraging the intermediate computational steps as external memory.

**The emergence of iteration heads.** First, [Cabannes et al. \(2024\)](#) identified the “iteration head” — an attention head that emerges during CoT. These heads explicitly focus on the model’s previously generated tokens to carry forward interim results. For example, in a loop counter task, an iteration head attends to the token “Step 4” to generate “Step 5”. This effectively allows the model to create a virtual *recurrent neural network* (RNN) where the hidden state is externalized as text. In another study of a *Llama-2* model ([Touvron et al., 2023](#)) solving multi-step ontology queries, [Dutta et al. \(2024\)](#) also identified early-layer attention heads that “move information along ontological relationships” in the contexts that are relevant to the current sub-problem. The emergence of iteration heads provides supporting evidence that CoT enables the model to internally utilize generated text as an external memory for sequential reasoning.

**Evidence of state maintenance and update.** Besides the access to external memory, studies show that LLMs with CoT can also maintain and update dynamic internal states to track the reasoning process. [Zhang et al. \(2025a\)](#) found that when using CoT for state-tracking tasks, LLMs embed an implicit finite state automaton in their hidden layers. Specific feed-forward neurons in later layers were found to correspond directly to discrete problem states, forming a circuit that reliably updates with each new reasoning step. This internal state representation is highly robust and works correctly even with noisy or incomplete CoT steps, suggesting

the model learns a resilient state-updating algorithm. By probing individual neurons of LLMs, [Rai and Yao \(2024\)](#) offered more granular evidence of state maintenance. They identified specific “reasoning neurons” in *Llama-2*’s feed-forward layers that activate to hold partial results, such as carried values during arithmetic. Their activation helps explain why including particular steps (*e.g.*, an explicit breakdown of a sum) in the CoT prompt is effective: they reliably trigger the neurons responsible for maintaining the intermediate state.

**Computational depth matters more than token semantics.** Notably, the internal process of sequential reasoning appears to persist even when the CoT rationale lacks semantic meaning. For example, [Pfau et al. \(2024\)](#) replaced the meaningful CoT text with filler tokens (*e.g.*, “...”). Surprisingly, models could still solve complex reasoning tasks simply by generating these dots. Similarly, [Goyal et al. \(2024\)](#) found that introducing a learnable “pause” token significantly boosts performance on tasks from QA to math. These findings suggest that the semantic content of reasoning steps may be secondary to the computational time they buy. The sheer act of generating extra tokens (regardless of their meaning) provides necessary computational depth; each token grants the model an additional forward pass through all its layers. This extra “think time” enables the model to implement complex reasoning algorithms that cannot be executed in a single pass. [Bharadwaj \(2024\)](#) reinforced this interpretation through a mechanistic study. They found that even when CoT steps are replaced by placeholders, the model’s deeper layers still encode the missing steps, which can be recovered to their correct semantic content via a logit lens probe.

**Parallelism and reasoning shortcuts.** Finally, although growing evidence reveals the sequential nature of CoT’s internal computation, other studies found that *LLMs often run multiple reasoning pathways in parallel during CoT*, so the model’s internal reasoning process is not strictly sequential. For example, [Dutta et al. \(2024\)](#) identified a “functional rift” where the model simultaneously tries to solve the problem directly from the question (“reasoning shortcuts”) while also following the step-by-step procedure, and these parallel approaches then converge in later layers. [Nikankin et al. \(2025\)](#) found that models perform arithmetic via many simple feature detectors rather than a single step-by-step algorithm. The models can still arrive at the correct

509 answer, even if they might make a mistake in an  
510 early step internally (Arcuschin et al., 2025). The  
511 above evidence reveals that CoT’s internal work-  
512 ings are more complicated than expected: it is a  
513 combination of sequential step-by-step reasoning,  
514 parallel associative shortcuts, and occasional after-  
515 the-fact rationalizations.

**Takeaway:** CoT activates a robust “reasoning mode” where models leverage generated tokens as external memory to execute *step-wise* internal computation, including retrieving intermediate results and updating internal states. However, this internal computation is not strictly sequential but a *parallel* process involving multiple pathways and shortcuts.

### 516 3.3 Why CoT Enhances Reasoning Abilities?

517 Explicit reasoning with CoT often solves complex  
518 tasks more accurately than implicit reasoning. Sev-  
519 eral reasons have been identified for why CoT can  
520 better elicit reasoning capabilities.  
521

522 **CoT augments computational expressiveness.**  
523 Recent theoretical studies demonstrate that CoT  
524 enhances transformers’ expressiveness and compu-  
525 tational capacity, enabling them to solve problems  
526 in higher complexity classes. A standard trans-  
527 former decoder without CoT performs constant-  
528 depth computation per token, limiting it to the com-  
529 plexity class  $TC^0$  (Merrill and Sabharwal, 2023a,b;  
530 Chiang et al., 2023). Such models theoretically  
531 cannot solve inherently serial problems because  
532 the required computation depth grows with input  
533 size, while the model’s depth is fixed. CoT breaks  
534 this limit. By feeding the output back into the input,  
535 CoT allows the transformer to simulate an RNN or  
536 a Turing Machine. The effective depth of the com-  
537 putation becomes proportional to the length of the  
538 generated chain. This elevates the transformer’s  
539 expressiveness to Polynomial Time (P) (Merrill  
540 and Sabharwal, 2024), making inherently serial or  
541 recursive computations solvable where they other-  
542 wise are not (Feng et al., 2023; Li et al., 2024b;  
543 Kim and Suzuki, 2025; Bavandpour et al., 2025).

544 **CoT introduces modularity that reduces sample**  
545 **complexity.** CoT decomposes complex tasks into  
546 granular, independent sub-problems. This mod-  
547 ularity provides an inductive bias that matches  
548 the structure of complex, multi-step problems, en-  
549 abling the model to master tasks with significantly

550 less data. Through both experimental and theoret-  
551 ical evidence, Li et al. (2023) demonstrated that  
552 CoT decouples in-context learning into a “filter-  
553 ing” phase and a “learning” phase that significantly  
554 reduces the sample complexity required to learn  
555 compositional structures like MLPs. Extending  
556 this learnability perspective, Yang et al. (2025a)  
557 demonstrated that CoT can render inherently “un-  
558 learnable” tasks efficiently learnable by reducing  
559 the sample complexity of the overall task to that  
560 of its hardest individual reasoning step. Wen et al.  
561 (2025) further identified that this efficiency stems  
562 from the *sparse sequential dependencies* among  
563 tokens. CoT induces interpretable, sparse attention  
564 patterns that enable polynomial sample complexity,  
565 whereas implicit reasoning requires exponentially  
566 many samples to disentangle dense dependencies.

**CoT enables more robust reasoning.** First, ev-  
567 idence shows that CoT *promotes robust general-*  
568 *ization* by encouraging models to learn general-  
569 izable solution patterns rather than overfitting to  
570 surface-level statistical shortcuts. For example, Yao  
571 et al. (2025a) demonstrated that CoT-trained mod-  
572 els induce a two-stage generalizing circuit that in-  
573 ternalizes the reasoning process, leading to strong  
574 OOD generalization even in the presence of train-  
575 ing noise. Complementing this, Li et al. (2025a)  
576 provided a theoretical guarantee for CoT general-  
577 ization, showing that CoT maintains high perfor-  
578 mance even when context examples are noisy or  
579 erroneous. Second, CoT helps *reduce the propaga-*  
580 *tion of errors* during reasoning. Gan et al. (2025)  
581 identified a “snowball error effect” in implicit rea-  
582 soning, where minor inaccuracies accumulate into  
583 significant failures. They demonstrated that CoT  
584 mitigates this by expanding the reasoning search  
585 space, which effectively lowers the probability of  
586 cumulative information loss and prevents errors  
587 from cascading through the reasoning chain.  
588

**Takeaway:** CoT enhances reasoning through three primary mechanisms: 1) It breaks the constant-depth limitation of the standard transformer, extending its effective computational depth. 2) It introduces modularity and inductive bias more aligned with multi-step reasoning, thus reducing the sample complexity required to learn complex tasks. 3) It facilitates robust generalization to OOD data and mitigates error propagation.

### 3.4 Does Chain-of-Thought Equate to Explainability?

Explicit reasoning appears to provide transparency, leading users to assume that CoT explanations accurately reveal how the model arrived at an answer. However, substantial evidence indicates that CoT outputs often do not faithfully reflect the model’s actual decision-making process (Lanham et al., 2023; Turpin et al., 2023; Chen et al., 2025c; Barez et al., 2025), a phenomenon referred to as the *unfaithfulness* of CoT reasoning.

**Evidence of CoT unfaithfulness.** Recent studies reveal that CoT frequently functions as *post-hoc rationalization* rather than the causal driver of predictions (Kudo et al., 2024; Arcuschin et al., 2025; Lewis-Lim et al., 2025). For instance, Turpin et al. (2023) demonstrated that models often alter their predictions based on spurious cues, such as the reordering of multiple-choice options. In such cases, the models still tend to confabulate logical-sounding CoT rationales that hide the actual spurious cause of their decision. Similarly, when correct answers are injected as hints, models often invent spurious derivations to support the injected answer without acknowledging the hint’s influence (Chen et al., 2025c). Furthermore, mechanistic analyses uncovered “silent error corrections”, where models internally correct mistakes without updating the CoT rationale (Arcuschin et al., 2025). Unfaithfulness is also evident in *sycophancy*, where models prioritize agreement with user beliefs over truthfulness. Even when models possess the correct internal knowledge, they frequently concede to incorrect user premises and generate plausible rationales to justify these compliant responses (Sharma et al., 2024). Collectively, these findings highlight a fundamental disconnect between verbalized rationales and internal computations, challenging the premise that CoT equates to explainability.

**Mechanistic understanding of CoT unfaithfulness.** Recent mechanistic analyses attribute this unfaithfulness to a fundamental mismatch between the distributed, parallel nature of transformer computation and the sequential nature of explicit reasoning. As discussed in Section 3.2, many works have revealed the *distributed nature* of LLMs’ internal reasoning; transformer-based LLMs frequently employ multiple redundant computational pathways to process information, *e.g.*, simultaneously leveraging memorization, heuristics, and al-

gorithmic circuits (McGrath et al., 2023; Dutta et al., 2024; Nikankin et al., 2025). Consequently, CoT only acts as a “lossy projection” of high-dimensional internal states, often capturing only a fraction of the model’s actual decision process (Dutta et al., 2024). Because computation is highly distributed, a single CoT rationale can capture at most one of many simultaneous causal pathways. As a result, CoTs typically omit influential factors and serve only as partial, post-hoc rationalisations of the model’s underlying distributed, superposed computation (Barez et al., 2025). This architectural dissonance makes unfaithfulness difficult to mitigate. Tanneru et al. (2024) demonstrated that even when training objectives explicitly penalize inconsistency, models still revert to plausible-but-not-causal explanations on complex tasks, highlighting the inherent difficulty in eliciting faithful CoT reasoning from LLMs.

**Takeaway:** While chain-of-thought offers the appearance of transparency, it does not equate to faithful explainability. CoT often functions as post-hoc rationalization rather than a true reflection of the model’s internal processing. Mechanistically, this unfaithfulness stems from a structural mismatch between the distributed, parallel computation of transformers and the sequential nature of explicit reasoning.

## 4 Future Research Directions

In this survey, we aim to provide a comprehensive overview of the mechanisms underlying multi-step reasoning in LLMs, around two fundamental reasoning paradigms: implicit reasoning and explicit reasoning. We systematically synthesize existing mechanistic studies through a framework of seven interconnected research questions.

Despite significant progress in opening the black box, critical challenges remain. As illustrated in Figure 1, we identify five strategic directions that are essential for the future of mechanistic understanding. *In Appendix A, we discuss these directions in detail*, including the need for rigorous causal analysis in real-world settings, bridging the faithfulness gap of explicit CoT, and developing white-box evaluation metrics. We hope this survey bridges the gap between empirical capabilities and mechanistic understanding, offering a unified perspective that guides the community toward building more reliable and controllable reasoning models.

## 681 Limitations

682 This survey concentrates strictly on the mechanistic  
683 understanding of multi-step reasoning within  
684 transformer-based LLMs. Consequently, we do not  
685 cover other aspects of reasoning, such as proba-  
686 bilistic inference, creative planning, or common-  
687 sense reasoning, which may operate under different  
688 mechanistic principles. Additionally, our scope is  
689 limited to the current paradigm of text-based trans-  
690 formers; we do not extensively address reasoning  
691 mechanisms in Multimodal LLMs (MLLMs), alter-  
692 native architectures like Diffusion Language Mod-  
693 els (DLMs), or neural networks that predate the  
694 modern era of large language models.

## 695 References

696 Roman Abramov, Felix Steinbauer, and Gjergji Kasneci.  
697 2025. [Grokking in the wild: Data augmentation for  
698 real-world multi-hop reasoning with transformers](#). In  
699 *Proceedings of the 42th International Conference on  
700 Machine Learning (ICML)*. OpenReview.net.

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

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

711 Alireza Amiri Bavandpour, Xinting Huang, Mark Rofin,  
712 and Michael Hahn. 2025. [Lower bounds for chain-  
713 of-thought reasoning in hard-attention transformers](#).  
714 In *Proceedings of the 42th International Conference  
715 on Machine Learning (ICML)*. OpenReview.net.

716 Aryasomayajula Ram Bharadwaj. 2024. [Understanding  
717 hidden computations in chain-of-thought reasoning](#).  
718 *CoRR*, abs/2412.04537.

719 Eden Biran, Daniela Gottesman, Sohee Yang, Mor Geva,  
720 and Amir Globerson. 2024. [Hopping too late: Ex-  
721 ploring the limitations of large language models on  
722 multi-hop queries](#). In *Proceedings of the 2024 Con-  
723 ference on Empirical Methods in Natural Language  
724 Processing (EMNLP)*, pages 14113–14130. Associa-  
725 tion for Computational Linguistics.

726 Jannik Brinkmann, Abhay Sheshadri, Victor Levoso,  
727 Paul Swoboda, and Christian Bartelt. 2024a. [A mech-  
728 anistic analysis of a transformer trained on a symbolic  
729 multi-step reasoning task](#). In *Findings of the Associa-  
730 tion for Computational Linguistics: ACL 2024*, pages  
731 4082–4102.

Jannik Brinkmann, Abhay Sheshadri, Victor Levoso,  
Paul Swoboda, and Christian Bartelt. 2024b. [A mech-  
anistic analysis of a transformer trained on a symbolic  
multi-step reasoning task](#). In *Findings of the Associa-  
tion for Computational Linguistics: ACL 2024*, pages  
4082–4102. Association for Computational Linguistics.

Vivien Cabannes, Charles Arnal, Wassim Bouaziz,  
Xingyu Yang, François Charton, and Julia Kempe.  
2024. [Iteration head: A mechanistic study of chain-  
of-thought](#). In *Proceedings of the 2024 Annual Con-  
ference on Neural Information Processing Systems  
(NeurIPS)*.

Yixin Cao, Jiahao Ying, Yaoning Wang, Xipeng Qiu,  
Xuanjing Huang, and Yugang Jiang. 2025. [Model  
utility law: Evaluating llms beyond performance  
through mechanism interpretable metric](#). *CoRR*,  
abs/2504.07440.

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

Xinghao Chen, Anhao Zhao, Heming Xia, Xuan Lu,  
Hanlin Wang, Yanjun Chen, Wei Zhang, Jian Wang,  
Wenjie Li, and Xiaoyu Shen. 2025b. [Reasoning be-  
yond language: A comprehensive survey on latent  
chain-of-thought reasoning](#). *CoRR*, abs/2505.16782.

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. 2025c. [Rea-  
soning models don’t always say what they think](#).  
*CoRR*, abs/2505.05410.

David Chiang, Peter Cholak, and Anand Pillay. 2023. [Tighter bounds on the expressivity of transformer  
encoders](#). In *Proceedings of the 40th International  
Conference on Machine Learning (ICML)*, volume  
202 of *Proceedings of Machine Learning Research*,  
pages 5544–5562.

Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang  
Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu,  
Bing Qin, and Ting Liu. 2024a. [Navigate through  
enigmatic labyrinth A survey of chain of thought rea-  
soning: Advances, frontiers and future](#). In *Proceed-  
ings of the 61th Annual Meeting of the Association for  
Computational Linguistics (ACL)*, pages 1173–1203.

Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang  
Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu,  
Bing Qin, and Ting Liu. 2024b. [Navigate through  
enigmatic labyrinth A survey of chain of thought  
reasoning: Advances, frontiers and future](#). In *Pro-  
ceedings of the 62nd Annual Meeting of the Asso-  
ciation for Computational Linguistics (ACL)*, pages  
1173–1203.

789	DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	846
790	Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	847
791	Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,	2021. <a href="#">Measuring massive multitask language under-</a>	848
792	Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-	<a href="#">standing</a> . In <i>Proceedings of the 9th International</i>	849
793	hong Shao, Zhuoshu Li, Ziyi Gao, and 81 others.	<i>Conference on Learning Representations (ICLR)</i> .	850
794	2025. <a href="#">Deepseek-r1: Incentivizing reasoning capa-</a>		
795	<a href="#">bility in llms via reinforcement learning</a> . <i>CoRR</i> ,	Yifan Hou, Jiaoda Li, Yu Fei, Alessandro Stolfo,	851
796	abs/2501.12948.	Wangchunshu Zhou, Guangtao Zeng, Antoine Bosse-	852
797	Subhabrata Dutta, Joykirat Singh, Soumen Chakrabarti,	lut, and Mrinmaya Sachan. 2023a. <a href="#">Towards a mech-</a>	853
798	and Tanmoy Chakraborty. 2024. <a href="#">How to think step-</a>	<a href="#">anistic interpretation of multi-step reasoning capa-</a>	854
799	<a href="#">by-step: A mechanistic understanding of chain-of-</a>	<a href="#">bilities of language models</a> . In <i>Proceedings of the</i>	855
800	<a href="#">thought reasoning</a> . <i>Transactions on Machine Learn-</i>	<i>2023 Conference on Empirical Methods in Natural</i>	856
801	<i>ing Research (TMLR)</i> , 2024.	<i>Language Processing (EMNLP)</i> , pages 4902–4919.	857
802	Yanai Elazar, Akshita Bhagia, Ian Magnusson, Abhi-	Yifan Hou, Jiaoda Li, Yu Fei, Alessandro Stolfo,	858
803	lasha Ravichander, Dustin Schwenk, Alane Suhr,	Wangchunshu Zhou, Guangtao Zeng, Antoine Bosse-	859
804	Evan Pete Walsh, Dirk Groeneveld, Luca Soldaini,	lut, and Mrinmaya Sachan. 2023b. <a href="#">Towards a mech-</a>	860
805	Sameer Singh, Hannaneh Hajishirzi, Noah A. Smith,	<a href="#">anistic interpretation of multi-step reasoning capa-</a>	861
806	and Jesse Dodge. 2024. <a href="#">What’s in my big data?</a> In	<a href="#">bilities of language models</a> . In <i>Proceedings of the</i>	862
807	<i>Proceedings of the 12th International Conference on</i>	<i>2023 Conference on Empirical Methods in Natural</i>	863
808	<i>Learning Representations (ICLR)</i> . OpenReview.net.	<i>Language Processing (EMNLP)</i> , pages 4902–4919.	864
809	Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye,	Association for Computational Linguistics.	865
810	Di He, and Liwei Wang. 2023. <a href="#">Towards revealing</a>	Jie Huang and Kevin Chen-Chuan Chang. 2023. <a href="#">To-</a>	866
811	<a href="#">the mystery behind chain of thought: A theoretical</a>	<a href="#">wards reasoning in large language models: A survey</a> .	867
812	<a href="#">perspective</a> . In <i>Proceedings of the 2023 Annual Con-</i>	In <i>Findings of the Association for Computational</i>	868
813	<i>ference on Neural Information Processing Systems</i>	<i>Linguistics: ACL 2023</i> , pages 1049–1065.	869
814	<i>(NeurIPS)</i> .	Mingyu Jin, Qinkai Yu, Dong Shu, Haiyan Zhao,	870
815	Zeyu Gan, Yun Liao, and Yong Liu. 2025. <a href="#">Rethinking</a>	Wenyue Hua, Yanda Meng, Yongfeng Zhang, and	871
816	<a href="#">external slow-thinking: From snowball errors to prob-</a>	Mengnan Du. 2024. <a href="#">The impact of reasoning step</a>	872
817	<a href="#">ability of correct reasoning</a> . In <i>Proceedings of the</i>	<a href="#">length on large language models</a> . In <i>Findings of</i>	873
818	<i>42th International Conference on Machine Learning</i>	<i>the Association for Computational Linguistics: ACL</i>	874
819	<i>(ICML)</i> . OpenReview.net.	<i>2024</i> , pages 1830–1842. Association for Computa-	875
820	Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lu-	tional Linguistics.	876
821	cas Dixon, and Mor Geva. 2024. <a href="#">Patchscopes: A</a>	Tianjie Ju, Yijin Chen, Xinwei Yuan, Zhuosheng Zhang,	877
822	<a href="#">unifying framework for inspecting hidden represen-</a>	Wei Du, Yubin Zheng, and Gongshen Liu. 2024. <a href="#">In-</a>	878
823	<a href="#">tations of language models</a> . In <i>Proceedings of the</i>	<a href="#">vestigating multi-hop factual shortcuts in knowledge</a>	879
824	<i>41th International Conference on Machine Learning</i>	<a href="#">editing of large language models</a> . In <i>Proceedings</i>	880
825	<i>(ICML)</i> . OpenReview.net.	<i>of the 62nd Annual Meeting of the Association for</i>	881
826	Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Kr-	<i>Computational Linguistics (ACL)</i> , pages 8987–9001.	882
827	ishna Menon, Sanjiv Kumar, and Vaishnavh Nagara-	Association for Computational Linguistics.	883
828	jan. 2024. <a href="#">Think before you speak: Training lan-</a>	Cheongwoong Kang and Jaesik Choi. 2023. <a href="#">Impact</a>	884
829	<a href="#">guage models with pause tokens</a> . In <i>Proceedings</i>	<a href="#">of co-occurrence on factual knowledge of large lan-</a>	885
830	<i>of the 12th International Conference on Learning</i>	<a href="#">guage models</a> . In <i>Findings of the Association for</i>	886
831	<i>Representations (ICLR)</i> . OpenReview.net.	<i>Computational Linguistics: EMNLP 2023</i> , pages	887
832	Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai	7721–7735. Association for Computational Linguis-	888
833	Dong, Wentao Zhang, Guanting Chen, Xiao Bi,	tics.	889
834	Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wen-	Juno Kim and Taiji Suzuki. 2025. <a href="#">Transformers prov-</a>	890
835	feng Liang. 2024. <a href="#">Deepseek-coder: When the large</a>	<a href="#">ably solve parity efficiently with chain of thought</a> . In	891
836	<a href="#">language model meets programming - the rise of code</a>	<i>Proceedings of the 13th International Conference on</i>	892
837	<a href="#">intelligence</a> . <i>CoRR</i> , abs/2401.14196.	<i>Learning Representations (ICLR)</i> . OpenReview.net.	893
838	Tianyu Guo, Hanlin Zhu, Ruiqi Zhang, Jiantao Jiao,	Keito Kudo, Yoichi Aoki, Tatsuki Kuribayashi, Shusaku	894
839	Song Mei, Michael I. Jordan, and Stuart Russell.	Sone, Masaya Taniguchi, Ana Brassard, Keisuke Sak-	895
840	2025. <a href="#">How do llms perform two-hop reasoning in</a>	aguchi, and Kentaro Inui. 2024. <a href="#">Think-to-talk or</a>	896
841	<a href="#">context?</a> <i>CoRR</i> , abs/2502.13913.	<a href="#">talk-to-think? when llms come up with an answer in</a>	897
842	Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li,	<a href="#">multi-step reasoning</a> . <i>CoRR</i> , abs/2412.01113.	898
843	Zhiting Hu, Jason Weston, and Yuandong Tian. 2024.	Tamera Lanham, Anna Chen, Ansh Radhakrishnan,	899
844	<a href="#">Training large language models to reason in a contin-</a>	Benoit Steiner, Carson Denison, Danny Hernan-	900
845	<a href="#">uous latent space</a> . <i>CoRR</i> , abs/2412.06769.	dez, Dustin Li, Esin Durmus, Evan Hubinger, Jack-	901
		son Kernion, Kamile Lukosiute, Karina Nguyen,	902





1129	Xiangqi Wang, Yue Huang, Yujun Zhou, Xiaonan Luo, Kehan Guo, and Xiangliang Zhang. 2025a. <a href="#">Causally-enhanced reinforcement policy optimization</a> . <i>CoRR</i> , abs/2509.23095.	1184
1130		1185
1131		1186
1132		1187
1133	Yiming Wang, Pei Zhang, Baosong Yang, Derek F. Wong, and Rui Wang. 2025b. <a href="#">Latent space chain-of-embedding enables output-free LLM self-evaluation</a> . In <i>Proceedings of the 13th International Conference on Learning Representations (ICLR)</i> . OpenReview.net.	1188
1134		1189
1135		1190
1136		1191
1137		1192
1138		1193
1139	Zijian Wang, Yanxiang Ma, and Chang Xu. 2025c. <a href="#">Eliciting chain-of-thought in base llms via gradient-based representation optimization</a> . <i>Preprint</i> , arXiv:2511.19131.	1194
1140		1195
1141		1196
1142		1197
1143	Jan Wehner, Sahar Abdelnabi, Daniel Tan, David Krueger, and Mario Fritz. 2025. <a href="#">Taxonomy, opportunities, and challenges of representation engineering for large language models</a> . <i>Transactions on Machine Learning Research (TMLR)</i> , 2025.	1198
1144		1199
1145		1200
1146		1201
1147		1202
1148		1203
1149	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022a. <a href="#">Emergent abilities of large language models</a> . <i>Transactions on Machine Learning Research (TMLR)</i> , 2022.	1204
1150		1205
1151		1206
1152		1207
1153		1208
1154		1209
1155	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022b. <a href="#">Chain-of-thought prompting elicits reasoning in large language models</a> . In <i>Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)</i> .	1210
1156		1211
1157		1212
1158		1213
1159		1214
1160		1215
1161	Kaiyue Wen, Huaqing Zhang, Hongzhou Lin, and Jingzhao Zhang. 2025. <a href="#">From sparse dependence to sparse attention: Unveiling how chain-of-thought enhances transformer sample efficiency</a> . In <i>Proceedings of the 13th International Conference on Learning Representations (ICLR)</i> . OpenReview.net.	1216
1162		1217
1163		1218
1164		1219
1165		1220
1166		1221
1167		1222
1168	Skyler Wu, Eric Meng Shen, Charumathi Badrinath, Jiaqi Ma, and Himabindu Lakkaraju. 2023. <a href="#">Analyzing chain-of-thought prompting in large language models via gradient-based feature attributions</a> . <i>CoRR</i> , abs/2307.13339.	1223
1169		1224
1170		1225
1171		1226
1172	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 22 others. 2024a. <a href="#">Qwen2.5 technical report</a> . <i>CoRR</i> , abs/2412.15115.	1227
1173		1228
1174		1229
1175		1230
1176		1231
1177		1232
1178		1233
1179	Chenxiao Yang, Zhiyuan Li, and David Wipf. 2025a. <a href="#">Chain-of-thought provably enables learning the (otherwise) unlearnable</a> . In <i>Proceedings of the 13th International Conference on Learning Representations (ICLR)</i> . OpenReview.net.	1234
1180		1235
1181		1236
1182		1237
1183		1238
	Sohee Yang, Elena Gribovskaya, Nora Kassner, Mor Geva, and Sebastian Riedel. 2024b. <a href="#">Do large language models latently perform multi-hop reasoning?</a> In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 10210–10229. Association for Computational Linguistics.	
	Sohee Yang, Nora Kassner, Elena Gribovskaya, Sebastian Riedel, and Mor Geva. 2025b. <a href="#">Do large language models perform latent multi-hop reasoning without exploiting shortcuts?</a> In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 3971–3992.	
	Sohee Yang, Nora Kassner, Elena Gribovskaya, Sebastian Riedel, and Mor Geva. 2025c. <a href="#">Do large language models perform latent multi-hop reasoning without exploiting shortcuts?</a> In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 3971–3992. Association for Computational Linguistics.	
	Zhipeng Yang, Junzhuo Li, Siyu Xia, and Xuming Hu. 2025d. <a href="#">Internal chain-of-thought: Empirical evidence for layer-wise subtask scheduling in llms</a> . <i>CoRR</i> , abs/2505.14530.	
	Xinhao Yao, Ruifeng Ren, Yun Liao, and Yong Liu. 2025a. <a href="#">Unveiling the mechanisms of explicit cot training: How chain-of-thought enhances reasoning generalization</a> . <i>CoRR</i> , abs/2502.04667.	
	Yuekun Yao, Yupei Du, Dawei Zhu, Michael Hahn, and Alexander Koller. 2025b. <a href="#">Language models can learn implicit multi-hop reasoning, but only if they have lots of training data</a> . <i>CoRR</i> , abs/2505.17923.	
	Jiaran Ye, Zijun Yao, Zhidian Huang, Liangming Pan, Jinxin Liu, Yushi Bai, Amy Xin, Liu Weichuan, Xiaoyin Che, Lei Hou, and Juanzi Li. 2025. <a href="#">How does transformer learn implicit reasoning?</a> In <i>Proceedings of the 2025 Annual Conference on Neural Information Processing Systems (NeurIPS 2025)</i> , San Diego, USA.	
	Xi Ye and Greg Durrett. 2022. <a href="#">The unreliability of explanations in few-shot prompting for textual reasoning</a> . In <i>Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)</i> .	
	Yijiong Yu. 2025. <a href="#">Do llms really think step-by-step in implicit reasoning?</a> <i>CoRR</i> , abs/2411.15862.	
	Zeping Yu, Yonatan Belinkov, and Sophia Ananiadou. 2025. <a href="#">Back attention: Understanding and enhancing multi-hop reasoning in large language models</a> . <i>CoRR</i> , abs/2502.10835.	
	Jason Zhang and Scott Viteri. 2024. <a href="#">Uncovering latent chain of thought vectors in language models</a> . <i>CoRR</i> , abs/2409.14026.	
	Yifan Zhang, Wenyu Du, Dongming Jin, Jie Fu, and Zhi Jin. 2025a. <a href="#">Finite state automata inside transformers with chain-of-thought: A mechanistic study on</a>	

1239 [state tracking](#). In *Proceedings of the 63rd Annual*  
1240 *Meeting of the Association for Computational Lin-*  
1241 *guistics (ACL)*, pages 13603–13621. Association for  
1242 Computational Linguistics.

1243 Yuyi Zhang, Boyu Tang, Tianjie Ju, Sufeng Duan,  
1244 and Gongshen Liu. 2025b. [Do latent tokens](#)  
1245 [think? a causal and adversarial analysis of chain-](#)  
1246 [of-continuous-thought](#). *CoRR*, abs/2512.21711.

1247 Zhongwang Zhang, Pengxiao Lin, Zhiwei Wang,  
1248 Yaoyu Zhang, and Zhi-Qin John Xu. 2025c. [Com-](#)  
1249 [plexity control facilitates reasoning-based compo-](#)  
1250 [sitional generalization in transformers](#). *CoRR*,  
1251 abs/2501.08537.

## 1252 A Future Research Directions

1253 As illustrated in Figure 1, we identify five strate-  
1254 gic directions that are essential for the future of  
1255 mechanistic understanding.

### 1256 **Rigorous causal analysis in real-world settings.**

1257 A fundamental challenge in current mechanistic  
1258 research is *the disparity between idealized experi-*  
1259 *mental settings and the complexities of real-world*  
1260 *reasoning*. First, the reliance on toy models and  
1261 synthetic data limits the generalizability of current  
1262 findings. For example, while the “grokking” phe-  
1263 nomenon has been identified as a potential pathway  
1264 for the emergence of implicit multi-hop reasoning,  
1265 most empirical evidence is derived from toy mod-  
1266 els trained from scratch on synthetic tasks (§ 2.2).  
1267 Consequently, it remains an open question whether  
1268 the phase transitions observed in these controlled  
1269 environments truly govern the development of rea-  
1270 soning capabilities in foundation models trained on  
1271 large-scale, naturalistic corpora.

1272 Second, the field should move beyond correla-  
1273 tional analysis, which only proves information pres-  
1274 ence, to *rigorous causal verification* within these  
1275 complex settings. Unlike clean synthetic environ-  
1276 ments, real-world data is ubiquitous with spuri-  
1277 ous cues, making it difficult to distinguish genuine  
1278 reasoning circuits from robust shortcut heuristics  
1279 (§ 2.3). Therefore, causal interventions are crucial  
1280 for proving that identified internal representations  
1281 are truly drivers of correct inference in the wild.  
1282 This understanding should ideally translate into *ro-*  
1283 *bust training-time interventions* that penalize such  
1284 shortcuts, forcing models to learn generalizable  
1285 algorithms despite the noisy data distribution. Ulti-  
1286 mately, future work must aim to synthesize these  
1287 insights into a unified theoretical framework that  
1288 explains how diverse components, from memoriza-

1289 tion circuits to reasoning heads, interact within the  
1290 massive scale of foundation models.

### 1291 **Bridging the faithfulness gap of explicit CoT**

1292 **reasoning.** As discussed in § 3.4, a critical bottle-  
1293 neck in current LLMs is the “functional rift” (Dutta  
1294 et al., 2024) between the model’s internal, par-  
1295 allel processing and its sequential, explicit CoT  
1296 reasoning. This structural mismatch forces mod-  
1297 els to compress high-dimensional, distributed la-  
1298 tent states into a low-bandwidth stream of dis-  
1299 crete tokens, often resulting in CoT that functions  
1300 as a post-hoc rationalization rather than a causal  
1301 driver. To address this, future research must ex-  
1302 plore *white-box alignment methods* that enforce  
1303 a causal link between implicit and explicit rea-  
1304 soning. Promising avenues include developing  
1305 training objectives that penalize discrepancies be-  
1306 tween the model’s hidden states (its true decision  
1307 process) and its generated rationale (Wang et al.,  
1308 2025a,c), imposing architectural constraints that  
1309 compel the model to rely solely on the generated  
1310 CoT for subsequent steps (Viteri et al., 2024), as  
1311 well as “self-explaining” dense internal represen-  
1312 tations into faithful natural language steps (Sengupta  
1313 and Rezik, 2025). Further exploration of these di-  
1314 rections is critical for aligning explicit outputs with  
1315 internal dynamics, ensuring CoT serves as a valid  
1316 window into the model’s computation.

### 1317 **Mechanistic understanding of Latent CoT rea-**

1318 **soning.** Beyond the dichotomy of implicit and  
1319 explicit CoT, an emerging paradigm is *latent CoT*  
1320 *reasoning* (Chen et al., 2025b; Li et al., 2025b),  
1321 where models are designed to simulate explicit rea-  
1322 soning trajectories entirely within hidden states.  
1323 Unlike standard implicit reasoning, which relies  
1324 on the fixed depth of a standard transformer, la-  
1325 tent CoT architectures often introduce additional  
1326 computational capacity via continuous “thought  
1327 tokens”, iterative refinement, or recurrent state up-  
1328 dates, frequently learning these behaviors by dis-  
1329 stilling explicit CoT data into latent representations.  
1330 This approach theoretically offers the best of both  
1331 worlds: it broadens the model’s expressive capacity  
1332 and computational depth while eliminating the re-  
1333 dundant decoding costs of natural language tokens.

1334 While various latent CoT architectures have been  
1335 proposed (Hao et al., 2024; Mitra et al., 2024;  
1336 Shen et al., 2025), mechanistic interpretability  
1337 has lagged significantly behind these innovations.  
1338 While a vast body of work has explored the latent  
1339 reasoning mechanisms of *standard transformers*

(§ 2.1), research into the internal dynamics of these novel latent CoT models remains limited (Zhang and Viteri, 2024; Wang et al., 2025b; Zhang et al., 2025b). Critical open questions remain: Does distilling explicit CoT truly force the model to internalize a sequential, step-by-step reasoning process, or does the model collapse the teacher’s rationale into high-dimensional statistical shortcuts? Therefore, gaining more mechanistic insights is crucial for designing next-generation latent CoT architectures and training objectives that effectively combine the interpretability of explicit reasoning with the efficiency of implicit computation.

**White-box evaluation metrics for LLM reasoning.** As we gain a deeper mechanistic understanding of multi-step reasoning, it should guide the development of evaluation protocols that go beyond simple end-task accuracy. Current black-box metrics (e.g., final accuracy) are increasingly insufficient, as models frequently arrive at correct answers via non-robust shortcuts, statistical heuristics, or “bag-of-words” processing (§ 2.3). To rigorously distinguish genuine reasoning from sophisticated pattern matching, the field requires “white-box” evaluation metrics that integrate model internals into the evaluation protocol. Pioneering efforts have begun to explore this direction. For example, Cao et al. (2025) introduced a mechanism-interpretable metric (MUI) that quantifies the “effort” required to solve a task, defined as the proportion of activated neurons or features. A truly capable model should achieve higher performance with lower effort. While this area remains underexplored, developing metrics that not only score the final output but also verify the presence of necessary internal computational signatures, such as the formation of bridge entities in intermediate layers (Yang et al., 2025c), is a crucial future trend. By defining reasoning not just as the correct outcome but as the execution of a verified internal process, we can prevent the overestimation of model capabilities and ensure that improvements on leaderboards reflect true algorithmic generalization.

**From mechanistic interpretation to model control.** While current research has successfully identified various reasoning circuits, such as iteration heads or deduction heads, most work remains observational. A major frontier for future study is the shift towards *pragmatic interpretability* (Nanda et al., 2025a,b), moving from passively explaining mechanisms to actively leveraging them for model

control and editing, a paradigm closely aligned with *Representation Engineering (RepE)* (Wehner et al., 2025). For example, if we can reliably identify the specific components responsible for multi-step logic, e.g., the state-maintenance neurons identified by (Rai and Yao, 2024), we can potentially intervene in real-time to correct reasoning errors or suppress shortcut neurons (Ju et al., 2024). Such interventions enable the development of “self-correcting” architectures that actively monitor internal states to detect and resolve failures like “silent errors” on the fly. Ultimately, this enables a transition from interpretability as a passive analysis tool to an active, foundational component for robust and safe reasoning systems.

1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403  
1404  
1405