

Navigate through Enigmatic Labyrinth

A Survey of Chain of Thought Reasoning: Advances, Frontiers and Future

Anonymous ACL submission

Abstract

Reasoning, a fundamental cognitive process in human intelligence, has garnered significant attention in the realm of artificial intelligence. Recent studies have found that chain-of-thought prompting significantly enhances LLM’s reasoning capabilities, which attracts widespread attention from both academia and industry. However, the field lacks a systematic survey. In this paper, we systematically investigate pertinent research, summarizing advanced methods from novel perspectives by meticulous taxonomy. Moreover, we delve into the current frontiers and delineate the challenges and future directions, thereby shedding light on future research. Furthermore, we engage in a discussion about open questions. We hope this paper serves as an introduction for beginners and fosters future research. Relevant resources have been made public available¹.

1 Introduction

In the realm of human cognition, reasoning stands as the linchpin, playing a vital role in the comprehension of the world and the formulation of decisions. As pre-training scales continue to expand (Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023a,b), language models exhibit growing capabilities (Wei et al., 2022a; Schaeffer et al., 2023; Zhou et al., 2023c), but challenges persist in the face of complex reasoning (Cobbe et al., 2021; Geva et al., 2021). Surprisingly, recent studies have found that guiding language models to reason step-by-step can enhance their ability to tackle intricate problems (Wei et al., 2022b; Jin and Lu, 2023), also known as chain-of-thought prompting (CoT). As depicted in Figure 1, the model progressively navigates its way out of the enigmatic labyrinth under the guidance of CoT prompting, finally arriving at the correct answer.

¹Resources are available at <https://github.com/>, updated periodically

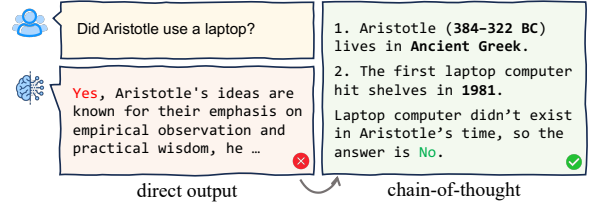


Figure 1: The model tackles complex problems step-by-step under the guidance of chain-of-thought prompting.

Thanks to the remarkable performance of CoT prompting, it has garnered widespread attention in both academia and industry, evolving into an independent research trajectory outside the realms of prompting engineering (Liu et al., 2023d; Qiao et al., 2023). Moreover, it has emerged as a crucial component in the landscape of AI autonomous agents (Wang et al., 2023h; Xi et al., 2023). However, these studies have yet to lack a systematic review and analysis. To fill this gap, we propose this work to conduct a comprehensive and detailed analysis of the XoT family. It’s worth noting that this paper explores the generalized chain-of-thought (XoT) from a broad perspective, with its core idea centered on reasoning step-by-step, progressively addressing complex problems.

Our contributions can be summarized as follows: (1) **First Survey**: This is the first comprehensive survey dedicated for XoT reasoning; (2) **Meticulous taxonomy**: We introduce a meticulous taxonomy (shown in Figure 2); (3) **Frontier and Future**: We discuss new frontiers, outline their challenges, and shed light on future research. (4) **Resources**: We make the resources publicly available to facilitate the research community.

Survey Organization We first give background and preliminary (§2); then present benchmarks (§3) and advanced methods (§4) from different perspectives. Furthermore, we discuss frontier research (§5) and outline challenges and future research directions (§6). Finally, we give a further discussion about open questions (§A.2).

2 Background and Preliminary

2.1 Background

In recent years, as model sizes increase (Brown et al., 2020; Scao et al., 2022; Touvron et al., 2023b; Zhao et al., 2023b), language models have emerged with numerous new capabilities, such as in-context learning (ICL) (Wei et al., 2022a; Brown et al., 2020) and chain-of-thought reasoning (Wei et al., 2022b). Accompanying this trend, pretrain with ICL has gradually supplanted pretrain with fine-tune, becoming the new paradigm in NLP (Qiu et al., 2020).

ICL integrates input-output demonstrations into prompts, enabling inference through few-shot learning. Through ICL, LLMs achieve competitive performance without fine-tuning but underperform in the face of complex reasoning tasks, while CoT prompting presents reasoning steps to LLMs, guiding them to solve complex problems progressively, thereby enhancing reasoning capabilities. Moreover, it exposes the LLM’s reasoning process to users, which offers interpretability.

2.2 Preliminary

In this section, we introduce the preliminary chain-of-thought reasoning with LLMs. Suppose there is a question \mathcal{Q} , a prompt \mathcal{T} and a probabilistic language model P_{LM} . The model takes the question and prompt as inputs to give the rationale \mathcal{R} and answer \mathcal{A} . We first consider in-context scenarios where the demonstrations do not contain reasoning chains. We need to maximize the likelihood of Answer \mathcal{A} , as shown in Equ. (1,2).

$$p(\mathcal{A} | \mathcal{T}, \mathcal{Q}) = \prod_{i=1}^{|\mathcal{A}|} p_{LM}(a_i | \mathcal{T}, \mathcal{Q}, a_{<i}) \quad (1)$$

$$\mathcal{T}_{ICL} = \{I, (x_1, y_1), \dots, (x_n, y_n)\} \quad (2)$$

In chain-of-thought reasoning scenario, where the demonstrations contain reasoning process, we need to maximize the likelihood of Answer \mathcal{A} and rationale \mathcal{R} , as shown in Equ. (3,4,5,6).

$$p(\mathcal{A} | \mathcal{T}, \mathcal{Q}) = p(\mathcal{A} | \mathcal{T}, \mathcal{Q}, \mathcal{R})p(\mathcal{R} | \mathcal{T}, \mathcal{Q}) \quad (3)$$

$$p(\mathcal{R} | \mathcal{T}, \mathcal{Q}) = \prod_{i=1}^{|\mathcal{R}|} p_{LM}(r_i | \mathcal{T}, \mathcal{Q}, r_{<i}) \quad (4)$$

$$p(\mathcal{A} | \mathcal{T}, \mathcal{Q}, \mathcal{R}) = \prod_{j=1}^{|\mathcal{A}|} p_{LM}(a_j | \mathcal{T}, \mathcal{Q}, \mathcal{R}, a_{<j}) \quad (5)$$

$$\mathcal{T}_{CoT} = \{I, (x_1, e_1, y_1), \dots, (x_n, e_n, y_n)\} \quad (6)$$

3 Benchmarks

Mathematical Reasoning Mathematical reasoning forms the foundation of human intelligence, playing a crucial role in problem-solving, decision-making, and world comprehension². It is commonly used to assess the general reasoning ability of language models (Patel et al., 2021; Cobbe et al., 2021; Hendrycks et al., 2021b; Mishra et al., 2022a).

Commonsense Reasoning Commonsense reasoning is essential for the interaction in daily life and the perception of the world, which assesses the world comprehension capacity of language models (Talmor et al., 2019, 2021; Geva et al., 2021).

Symbolic Reasoning Symbolic reasoning disentangles semantics and serves as a testbed for language models’ competence in simulating atomic operations (Wei et al., 2022b; Srivastava et al., 2022; Suzgun et al., 2023).

Logical Reasoning Logical reasoning is of paramount importance as it serves as the bedrock for rational thinking, robust problem-solving and interpretable decision-making (Liu et al., 2020; Yu et al., 2020; Tafjord et al., 2021; Han et al., 2022).

Multi-modal Reasoning Multimodal reasoning goes beyond the text, connecting human thought (text) with the natural world (vision, auditory, etc.) (Zellers et al., 2019; Park et al., 2020; Xiao et al., 2021; Lu et al., 2022).

4 Advanced Methods

In this section, we will discuss advanced XoT methods from three perspectives: construction approach (§4.1), structural variations (§4.2), and enhancement methods (§4.3). The taxonomy is shown in Figure 2.

4.1 Construction Approach

Based on the human effort required for model performing XoT reasoning, we divide the construction approaches into three categories: 1) Manual XoT, 2) Automatic XoT, and 3) Semi-automatic XoT.

4.1.1 Manual XoT

Wei et al. (2022b) first propose chain-of-thought prompting (Fewshot CoT) by manually annotating natural language form rationales to guide models in

²Please refer to Appendix for details of benchmarks.

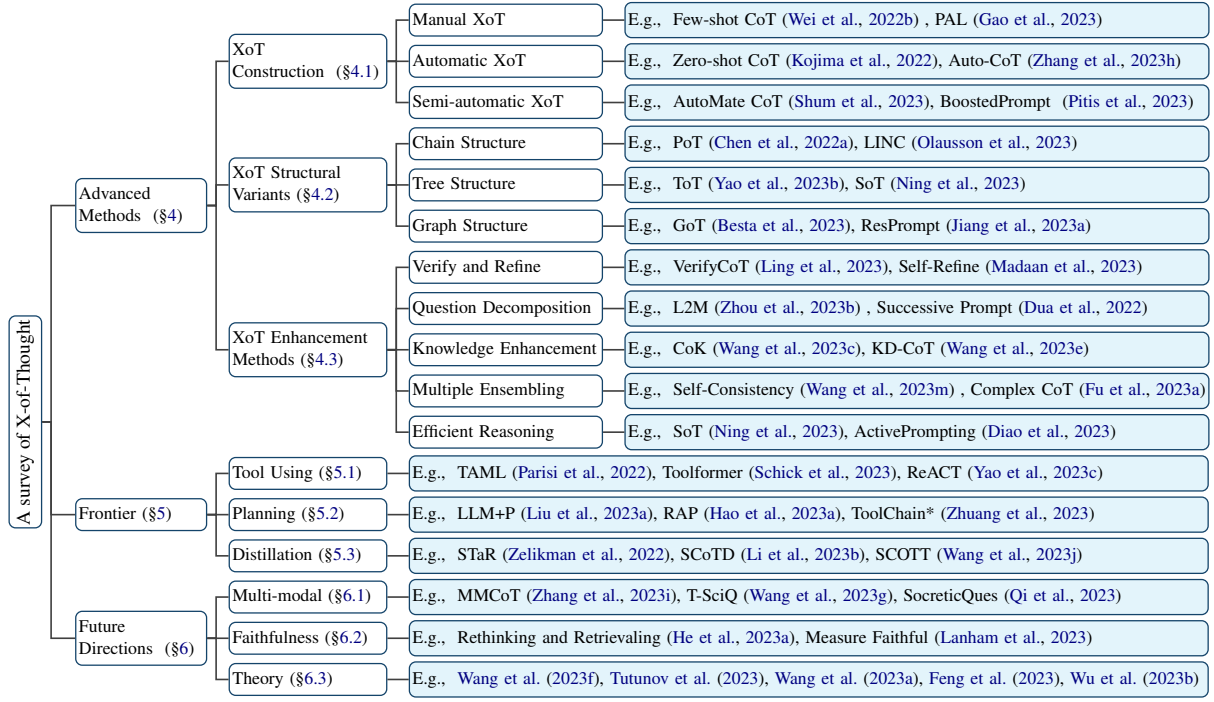


Figure 2: Taxonomy of Advanced Methods, Frontiers and Future Directions (Full version in Figure 8).

stepwise reasoning. To mitigate intermediate errors in reasoning, PAL (Gao et al., 2023), PoT (Chen et al., 2022a), MathPrompter (Imani et al., 2023) and NLEP (Zhang et al., 2023d) leverage rationales in programming language form, transforming problem-solving into program generation, and obtaining a deterministic answer through external program executor. Moreover, Fu et al. (2023a) discovers that using complex reasoning chains as demonstrations can further improve reasoning performance.

4.1.2 Automatic XoT

Some work designs specific instructions to stimulate CoT reasoning in a zero-shot manner, such as appending *Let's think step by step* after questions (Zeroshot CoT) (Kojima et al., 2022). There are also other types of instructions, including writing programs (Chen et al., 2022a), creating plans (Wang et al., 2023i), and generating task-related descriptions (Crispino et al., 2023), etc.

However, due to the lack of demonstration guidance, instruction-based methods are extremely unstable. Another route of work conducts few-shot reasoning based on automatically generated rationales (usually by Zeroshot CoT), which provides more stable reasoning. Such approaches focus on demonstration selection to boost reasoning. Zhang et al. (2023h) chooses diverse rationales through clustering, Zou et al. (2023) builds demonstrations

based on the question pattern, Wan et al. (2023) employs answer entropy as a metric for selection, and Xu et al. (2023) uses gibbs sampling to iteratively select samples.

4.1.3 Semi-automatic XoT

Building upon automatic methods rooted in few-shot learning, semi-automatic approaches incorporate a small number of human-annotated rationales to obtain supervision signals. They focus on bootstrapping to acquire high-quality rationales and selecting appropriate demonstrations to facilitate reasoning. Shao et al. (2023b) generates high-quality rationales through alternating forward and backward synthetic processes and Pitis et al. (2023) iteratively expands the examples when encountering challenging questions, which mitigates the issue of limited human supervision. On the other hand, some studies optimize demonstration selection. Shum et al. (2023) and Lu et al. (2023b) utilize policy gradient strategy to find examples, while Ye and Durrett (2023) applies two proxy metrics on development sets to yield demonstrations.

4.1.4 Pros and Cons of three Approaches

Manual XoT relies on high-quality rationale annotations, which result in superior performance. However, it encounters drawbacks such as high labor requirements and challenges in domain transfer. In contrast, Automatic XoT incurs no labor costs and allows for free domain transfer. How-

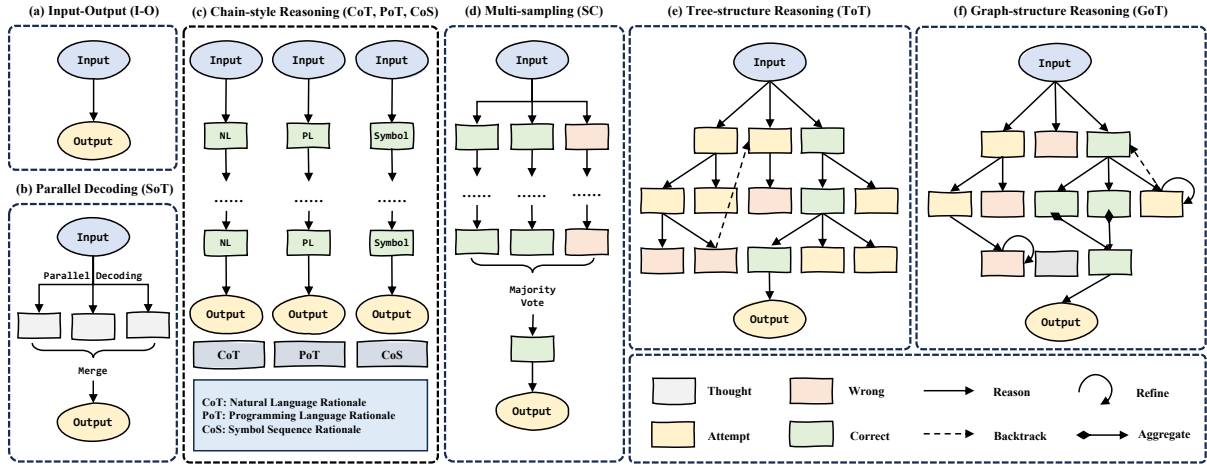


Figure 3: Structural variants emerging in the evolution of XoT. (a) standard I-O prompting (b) parallel-constrained tree structure variants (c) chain structure variants with distinct rationale descriptions (d) chain structure variants with multiple sampling (e) standard tree structure variants (f) standard graph structure variants.

ever, it is plagued by errors and instability due to a lack of supervised signals. Semi-automatic XoT strikes a subtle balance between the two, achieving a nuanced trade-off between performance and costs, making it more suitable for real-world applications.

4.2 XoT Structural Variants

The evolution of XoT has led to the development of multiple topological variants³. In this section, we will introduce variants of chain structure, tree structure and graph structure.

Chain Structure The descriptive form of rationales significantly influences reasoning execution. PAL (Gao et al., 2023) and PoT (Bi et al., 2023) use programming languages to depict the reasoning process, turning problem-solving into code generation. Similarly, formal logic description languages are also used to depict logical reasoning (Olausson et al., 2023; Pan et al., 2023; Ye et al., 2023a). They decouple the thought generation from execution, eliminating inconsistency errors. Additionally, algorithmic descriptions (Sel et al., 2023) can offer a high-level reasoning framework instead of addressing specifics, endowing the model with the ability for global thinking.

Tree Structure Chain structure inherently limits the scope of exploration. Through the incorporation of tree structures and search algorithms, models gain the capability to widely explore and backtrack during reasoning (Long, 2023; Yao et al., 2023b), as shown in Figure 3(e). Benefiting from

the exploration, tree variants have gained preliminary global planning capabilities towards the global optimum. Meanwhile, (Mo and Xin, 2023; Cao et al., 2023) introduce uncertainty measures based on Monte Carlo dropout and generation likelihood, respectively, thereby offering a more accurate evaluation of intermediate reasoning processes. To address complex problems, Yu et al. (2023b) uses a bottom-up approach by building an analogy sub-problems tree. In addition, Ning et al. (2023) accelerates reasoning by solving tree structure sub-problems in parallel. However, current methods are restricted by demands of explicit question decomposition and state transition, which leads to limitations in task generalization.

Graph Structure Graph structures introduce loops and N-to-1 connections, enabling improved modeling of subproblem aggregation and self-verification (Besta et al., 2023; Lei et al., 2023a), as illustrated in Figure 3(f). When confronted with complex problems, it demonstrates superior performance compared to tree variants, but faces similar challenges in task generalization. To address the generalization, Jiang et al. (2023a) establishes connections between reasoning steps in the prompts, thereby implicitly constructing a reasoning graph, which alleviates constraints imposed by complex topological structure.

The models’ capability progresses as the structure becomes more complex. Nevertheless, the generalization is limited by complex topological structures. The primary challenge for future research lies in extending methods based on these complex structures to universal domains.

³We consider XoT with chain structure and natural lang. rationales as vanilla CoT (the most primitive chain-of-thought).

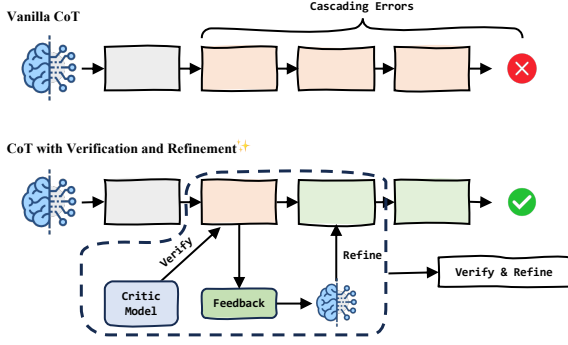


Figure 4: Verification and refinement rectify intermediate errors, which reduce cascading errors in reasoning.

4.3 XoT Enhancement Methods

In the following, we introduce enhanced XoT methods from five perspectives, including verify and refine (§4.3.1), question decomposition (§4.3.2), knowledge enhancement (§4.3.3), multiple ensembling (§4.3.4) and efficient reasoning (§4.3.5).

4.3.1 Verify and Refine

LLMs tend to be hallucinatory, which manifests as factual and faithful errors in reasoning (Huang et al., 2023c). Incorporating verification and refinement can be an effective strategy for mitigating the phenomena. In this section, we primarily focus on mitigating faithful errors, with a separate discussion of factual errors in the following knowledge enhancement section (§4.3.3).

LLMs can refine reasoning based on critics’ feedback. Paul et al. (2023) trains a small critic model to provide structured feedback, but the quality of the feedback is limited due to the model size. Madaan et al. (2023) employs feedback from itself for iterative self-refinement, Li et al. (2022c) uses finer-grained feedback at the step level, and Shinn et al. (2023) further enhances this approach by incorporating long and short-term memory to provide more concise suggestions. However, recent research suggests that LLMs may not address issues beyond their own capabilities (Kadavath et al., 2022; Yin et al., 2023), which raises doubt on the effectiveness of self-feedback (Huang et al., 2023b). To address this, some work incorporates external feedback (Gou et al., 2023a; Nathani et al., 2023) or performs secondary verification on the refinement (Shridhar et al., 2023).

On the other hand, logical reasoning structures are also well-suited for verification. Ling et al. (2023) devises a deductive reasoning form named Natural Program, which guarantees that the conclusion is derived from the designated premises.

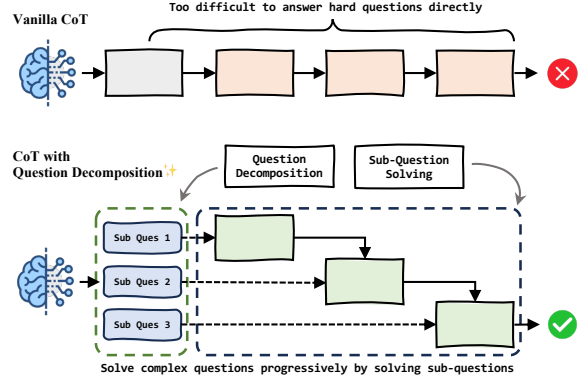


Figure 5: Question decomposition solves complex questions progressively by solving simple sub-questions.

Besides, backward (abductive) reasoning excels in detecting inconsistencies in reasoning. It reconstructs conditions or variables in the question based on the reasoning chain to discover inconsistencies, thereby refining the reasoning (Xue et al., 2023; Weng et al., 2022; Jiang et al., 2023b).

Reasoning with LLMs is prone to hallucinations, and feedback from intermediate steps plays a crucial role in refining the reasoning. However, the current acquisition of feedback signals still has many shortcomings, which necessitates further research.

4.3.2 Question Decomposition

The core idea of XoT is to solve questions step-by-step. However, vanilla CoT does not explicitly decompose questions, making it challenging to answer complex questions. To address this, certain approaches embrace the divide-and-conquer philosophy, overcoming intricate problems by tackling straightforward sub-problems.

L2M (Zhou et al., 2023b) initially breaks down the question into sub-questions in a top-down fashion. It then solves one sub-question at a time and leverages its solution to facilitate subsequent sub-questions. Dua et al. (2022) takes a similar approach to L2M, but it uses solutions from previous sub-questions to iteratively decompose questions. Khot et al. (2023) designs a modular task-sharing library that tailors more effective solutions to different classes of sub-questions. In multi-hop reasoning, iterative decomposition has become a common practice (Wang et al., 2022; Press et al., 2023; Trivedi et al., 2023). Additionally, some methods obtain a dedicated decomposer through supervised training rather than relying on the LLM itself (Li et al., 2023d; Junbing et al., 2023). However, when dealing with tabular reasoning, answering sub-questions may also pose a challenge, par-

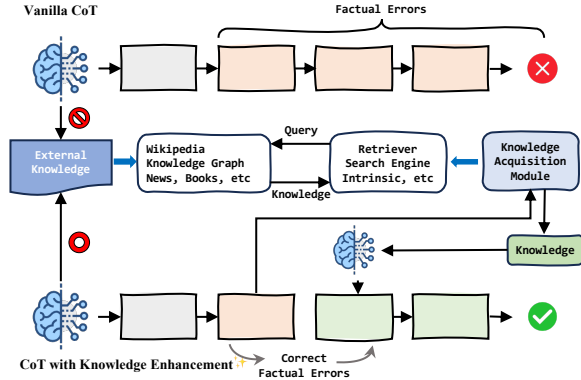


Figure 6: Introducing knowledge (external or internal) reduces factual errors in reasoning.

ticularly when handling large tables. To tackle this issue, certain approaches involve decomposing both the questions and tables simultaneously (Ye et al., 2023b; Cheng et al., 2023).

In addition to top-down decomposition, bottom-up sub-question aggregation is also a viable solution, with a smaller exploration space that leads to lower costs. Qi et al. (2023) employs Socratic questioning for recursive self-questioning to solve complex questions, while Zhang et al. (2023e), in a similar fashion, breaks tasks into small components and resolves them bottom-up.

4.3.3 Knowledge Enhancement

When dealing with knowledge-sensitive tasks, LLMs often make factual errors. Introducing external knowledge or mining the model’s internal knowledge can help alleviate this issue. Some methods explicitly utilize the model’s intrinsic knowledge. For example, Dhuliawala et al. (2023); Ji et al. (2023); Zheng et al. (2023c) prompt models to output its parametric knowledge, and then reason based on that intrinsic knowledge. Additionally, Zhang et al. (2023g) prompts the model to perform inductive reasoning on its intrinsic knowledge, deriving more general conclusions. Furthermore, Liu et al. (2023c) incorporates reinforcement learning to optimize based on model’s intrinsic knowledge. Meanwhile, Li and Qiu (2023) constructs an external memory bank using model’s reasoning chains and retrieves from it as needed.

External knowledge is often more reliable than parametric knowledge. Li et al. (2023d); Wang et al. (2023e) generates queries based on the question, utilizing a knowledge base as the external knowledge. Building upon this, Wang et al. (2023c) introduces a verification step for the retrieved knowl-

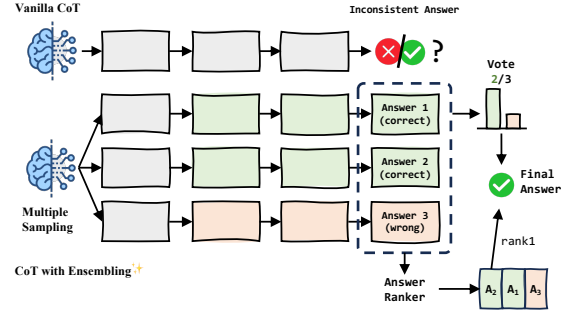


Figure 7: Voting and ranking reduce inconsistency by selecting final answers from multiple samplings.

edge, further ensuring knowledge accuracy. However, when confronted with multi-hop reasoning, direct retrieval using the question can be insufficient. Therefore, Press et al. (2023); Trivedi et al. (2023); Shao et al. (2023a); Yoran et al. (2023) decompose the question and iteratively use sub-question for more precise retrieval.

4.3.4 Multiple Ensembling

The sampling of generation introduces uncertainty, which, in turn, creates the possibility of improving performance through ensemble learning. Cobbe et al. (2021) trains a verifier to rank answers. SC (Wang et al., 2023m) performs majority voting based on answers across multiple samples, and Fu et al. (2023a) proposes a complexity-based voting strategy on top of SC. Widespread practical evidence indicates that ensemble is an effective way to improve performance. However, answer-based ensembling fails to consider intermediate steps. In response, Miao et al. (2023); Yoran et al. (2023); Khalifa et al. (2023) refines the ensemble at the step level. Yet another concern is the limited diversity offered by probability sampling. To overcome this limitation, Naik et al. (2023) uses different instructions, Liu et al. (2023e) ensembles various XoT variants, and Qin et al. (2023a) ensembles using multi-lingual reasoning chains. Furthermore, the multi-agent debate (MAD) framework can also be regarded as heterogeneous ensemblings (Liang et al., 2023; Du et al., 2023; Wang et al., 2023b).

4.3.5 Efficient Reasoning

LLMs are often inefficient in reasoning, such as high latency, substantial annotation costs, and elevated inference costs. To speed up reasoning, Ning et al. (2023) decomposes the questions in parallel and handles them simultaneously, Zhang et al. (2023b) generates a draft to skip intermediate layers during inference, and Leviathan et al. (2023)

introduces speculative decoding, which employs a smaller model for approximate inference. [Diao et al. \(2023\)](#) annotates high-uncertainty samples to reduce human costs, and [Aggarwal et al. \(2023\)](#) dynamically adjusts sampling frequency to reduce inference costs. Further research should focus on efficient reasoning to promote the widespread application of LLMs.

5 Frontiers Research

5.1 Tool Using

LLMs face challenges accessing news, performing calculations, and interacting with the environment. Previous work ([Parisi et al., 2022](#); [Schick et al., 2023](#); [Shen et al., 2023a](#)) grant LLM the ability to employ external tools, augmenting reasoning capabilities and assimilate external knowledge, enabling it to engage in calculation or multimodal interaction. However, the above approaches have limitations in facilitating multiple invocations of the tool and rectifying query errors. To tackle this problem, ReAct ([Yao et al., 2023c](#)) and Reflexion ([Shinn et al., 2023](#)) integrate the strengths of reasoning and action to complement each other. ART ([Paranjape et al., 2023](#)) uses a task library to select relevant tools and reasoning demonstrations.

These research studies focus on designing tools (or APIs) to enhance the capabilities of LLMs in various domains. Action facilitates interaction with external sources, such as knowledge bases and environments, enabling it to gather additional information. Simultaneously, XoT enables effective elicitation, tracking, and action refining.

5.2 Planning

LLMs cannot directly provide accurate responses for intricate problems, which requires planning to decompose them into sub-tasks and trace intermediate results. A plan can be described by code or definition language. AdaPlanner ([Sun et al., 2023](#)) generates Python code to control the agent and refines the plan iteratively based on feedback from execution. LLM+P ([Liu et al., 2023a](#)) and LLM+DP ([Dagan et al., 2023](#)) facilitate the Planning Domain Definition Language (PDDL) ([Gerevini, 2020](#)) to describe the planning procedure. PDDL assists in decomposing complex problems and utilizing specialized models for planning before converting the results into natural language. ISR-LLM ([Zhou et al., 2023d](#)) combines Self-Refine with PDDL to

achieve a better success rate in long-horizon sequential tasks.

Instead of pre-defined plans, many studies use search algorithms to plan and explore the action space dynamically. Tree-of-Thought ([Yao et al., 2023b](#)) decomposes the problem by deep-first or breadth-first search. Reasoning via Planning (RAP) ([Hao et al., 2023a](#)) and LATS ([Zhou et al., 2023a](#)) utilize LM-based Monte Carlo Tree Search for a more flexible planning procedure. Toolchain* ([Zhuang et al., 2023](#)) enables a more efficient exploring through heuristic A* search.

In summary, employing an LLM as a central controller, alongside tool usage and planning capabilities, constitutes a pathway toward realizing autonomous agents and, potentially, embodied intelligence in future research.

5.3 Distillation of Reasoning Capabilities

In low-resource scenarios such as edge computing, distillation offers a possibility for deploying LLMs. Besides, self-distillation is also a means of enhancing reasoning capabilities. [Huang et al. \(2023a\)](#) employs self-consistency to generate reasoning chains from unlabeled data, followed by fine-tuning, enhancing its generalized reasoning capabilities. [Zelikman et al. \(2022\)](#) improves LM’s reasoning capabilities via self-loop bootstrapping.

Though CoT is an emerging ability primarily observed in LLMs, it is limited in smaller models. [Magister et al. \(2023\)](#) demonstrates that fine-tuning T5 with reasoning chains generated by larger teacher models can substantially enhance task performance across diverse datasets. [Hsieh et al. \(2023b\)](#) generates rationales by prompting the language model to provide reasoning from the answer voted by self-consistency. [Ho et al. \(2023\)](#); [Li et al. \(2023b\)](#) finds that sampling multiple reasoning chains per instance is paramount for improving students’ capability. SCOTT ([Wang et al., 2023j](#)) utilizes contrastive decoding ([Li et al., 2022b](#); [O’Brien and Lewis, 2023](#)) and counterfactual reasoning objective to tackle the shortcut problem. DialCoT ([Han et al., 2023](#)) decomposes reasoning steps into a multi-round dialog and selects the correct path using the PPO algorithm.

These studies adopt a shared paradigm that distills smaller models with reasoning chains generated from larger models with superior reasoning capabilities. It’s notable that language models have intricate tradeoffs associated with multidimensional

capabilities, and distilling task-specific reasoning ability may adversely impact the general performance (Fu et al., 2023b).

6 Future Directions

While XoT has showcased remarkable performance on numerous tasks, there are still some challenges that necessitate further research.

6.1 Multi-modal Reasoning

Current XoT research mostly focuses on plain text. However, interacting with the real world necessitates multi-modal capabilities. To facilitate research, SciQA (Lu et al., 2022) and CURE (Chen et al., 2023b) are introduced to emphasize multi-modal CoT reasoning. Through fine-tuning with vision-language features, Zhang et al. (2023i); Wang et al. (2023g) endow models with multi-modal XoT capabilities, and Yao et al. (2023d,a) further incorporate graph structures to model multi-hop relationships. Other approaches convert images to captions and use LLM for prompt-based reasoning (Yang et al., 2023b; Zheng et al., 2023b). However, the limited capabilities of vision-language models constrain their performance in XoT reasoning (Alayrac et al., 2022; Li et al., 2023a; Peng et al., 2023).

Nevertheless, several critical challenges remain to be addressed in future research, which we summarize as follows: (1) **Visual-text interaction**: How can visual and textual features be effectively integrated, instead of relying solely on captions? (2) **Harnessing LLM**: How can we better apply LLM-based reasoning techniques to the multi-modal domain? (3) **Video Reasoning**: How to expand into video reasoning with complex temporal dependencies?

6.2 Faithful Reasoning

Extensive research indicates that LLMs often engage in unfaithful reasoning, such as factual errors and inconsistencies. To address factual errors, one common approach is retrieval augmentation (Trivedi et al., 2023; Zhao et al., 2023a), but it requires appropriate timing and retrieval accuracy. Compared to factual errors, inconsistencies are more difficult to identify. Common detection methods include logic-based (Jiang et al., 2023b; Xue et al., 2023; Ling et al., 2023), post-processing (He et al., 2023a; Lei et al., 2023b), and critic-based approaches (Madaan et al., 2023; Nathani et al.,

2023). Neural-symbolic reasoning (Chen et al., 2022a; Olausson et al., 2023) is a widely used approach for reducing inconsistencies, and question decomposition (Radhakrishnan et al., 2023) has also demonstrated its effectiveness to some degree. Furthermore, Zhang et al. (2023c); Lanham et al. (2023) investigate the factors influencing faithfulness from an empirical perspective.

The faithful reasoning encounters two significant challenges: (1) **Detection**: How can unfaithful reasoning be accurately identified? (2) **Correction**: How can precise feedback be generated to facilitate accurate correction?

6.3 Theoretical Perspective

The mechanism behind the CoT and ICL has not been clearly explained so far. Some studies empirically explore the roles of CoT and ICL in reasoning, offering practical insights (Wang et al., 2023a; Madaan and Yazdanbakhsh, 2022; Tang et al., 2023). Another line of work explores from a theoretical perspective. Li et al. (2023e); Feng et al. (2023); Merrill and Sabharwal (2023) investigate why CoT enhances reasoning abilities, while Wu et al. (2023b); Tutunov et al. (2023); Hou et al. (2023); Wang et al. (2023f) examine the mechanisms from a feature-based standpoint (information flow, attention). Additionally, there has been preliminary exploration of the emergence mechanism (Schaeffer et al., 2023; Zhou et al., 2023c).

At present, the exploration of CoT theories is still limited to the surface level. There are still open questions that require further in-depth investigation. (1) How does the **emergence capability** arise? (2) **In what way** does CoT enhance reasoning compared to standard ICL?

7 Conclusion

In this paper, we present a systematic survey of existing research on X-of-thought reasoning, offering a comprehensive review of the field. Specifically, we summarize and discuss the advanced methods from various perspectives. Additionally, we delve into the current frontiers, highlighting existing challenges, identifying potential research directions for the future, and discussing open questions⁴. This paper is the first systematic survey dedicated to XoT reasoning. We hope that this survey will facilitate further research in this area.

⁴Due to page limit, we leave related work, discussion in Appendix A, and benchmarks details in Appendix B

Limitations

This study provides the first comprehensive survey of generalized chain-of-thought (XoT) reasoning. Related work, benchmarks details and further discussion can be found in Appendix A,B.

We have made our best effort, but there may still be some limitations. Due to page limitations, we cannot provide every technical detail. We primarily gather studies from ACL*, NeurIPS, ICLR, ICML and arXiv, and there is a chance that we may have missed some important work published in other venues. In the benchmarks section, we primarily include widely used datasets, and more complete benchmarks can be found in Guo et al. (2023). As of now, there is no definitive conclusion on open questions. We will stay abreast of discussions within the research community, updating opinions and supplementing overlooked work in the future.

References

- Pranjal Aggarwal, Aman Madaan, Yiming Yang, and Mausam. 2023. [Let’s sample step by step: Adaptive-consistency for efficient reasoning with llms](#). *ArXiv preprint*, abs/2305.11860.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022. [Flamingo: a visual language model for few-shot learning](#). In *NeurIPS*.
- Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. [MathQA: Towards interpretable math word problem solving with operation-based formalisms](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics.
- Simran Arora, Avani Narayan, Mayee F. Chen, Laurel J. Orr, Neel Guha, Kush Bhatia, Ines Chami, and Christopher Ré. 2023. [Ask me anything: A simple strategy for prompting language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefer. 2023. [Graph of thoughts: Solving elaborate problems with large language models](#). *ArXiv preprint*, abs/2308.09687.
- Sumithra Bhakthavatsalam, Daniel Khashabi, Tushar Khot, Bhavana Dalvi Mishra, Kyle Richardson, Ashish Sabharwal, Carissa Schoenick, Oyvind Tafjord, and Peter Clark. 2021. [Think you have solved direct-answer question answering? try arc-da, the direct-answer AI2 reasoning challenge](#). *ArXiv preprint*, abs/2102.03315.
- Zhen Bi, Ningyu Zhang, Yinuo Jiang, Shumin Deng, Guozhou Zheng, and Huajun Chen. 2023. [When do program-of-thoughts work for reasoning?](#)
- Yonatan Bisk, Rowan Zellers, Ronan LeBras, Jianfeng Gao, and Yejin Choi. 2020. [PIQA: reasoning about physical commonsense in natural language](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7432–7439. AAAI Press.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott M. Lundberg, Harsha Nori, Hamid Palangi, Marco Túlio Ribeiro, and Yi Zhang. 2023. [Sparks of artificial general intelligence: Early experiments with GPT-4](#). *ArXiv preprint*, abs/2303.12712.
- Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. 2023. [Large language models as tool makers](#).
- Shulin Cao, Jiajie Zhang, Jiaxin Shi, Xin Lv, Zijun Yao, Qi Tian, Lei Hou, and Juanzi Li. 2023. [Probabilistic tree-of-thought reasoning for answering knowledge-intensive complex questions](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12541–12560, Singapore. Association for Computational Linguistics.

- Wenhu Chen. 2023. [Large language models are few\(1\)-shot table reasoners](#). In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1120–1130, Dubrovnik, Croatia. Association for Computational Linguistics.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2022a. [Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks](#). *ArXiv preprint*, abs/2211.12588.
- Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. 2023a. [TheoremQA: A theorem-driven question answering dataset](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7889–7901, Singapore. Association for Computational Linguistics.
- Yangyi Chen, Karan Sikka, Michael Cogswell, Heng Ji, and Ajay Divakaran. 2023b. [Measuring and improving chain-of-thought reasoning in vision-language models](#). *ArXiv preprint*, abs/2309.04461.
- Zhipeng Chen, Kun Zhou, Beichen Zhang, Zheng Gong, Wayne Xin Zhao, and Ji-Rong Wen. 2023c. [Chatcot: Tool-augmented chain-of-thought reasoning on chat-based large language models](#).
- Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. 2021. [FinQA: A dataset of numerical reasoning over financial data](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3697–3711, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Zhiyu Chen, Shiyang Li, Charese Smiley, Zhiqiang Ma, Sameena Shah, and William Yang Wang. 2022b. [ConvFinQA: Exploring the chain of numerical reasoning in conversational finance question answering](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6279–6292, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Zhoujun Cheng, Tianbao Xie, Peng Shi, Chengzu Li, Rahul Nadkarni, Yushi Hu, Caiming Xiong, Dragomir Radev, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Tao Yu. 2023. [Binding language models in symbolic languages](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Haotian Wang, Ming Liu, and Bing Qin. 2023. [Timebench: A comprehensive evaluation of temporal reasoning abilities in large language models](#). *ArXiv preprint*, abs/2311.17667.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *ArXiv preprint*, abs/2110.14168.
- Nicholas Crispino, Kyle Montgomery, Fankun Zeng, Dawn Song, and Chenguang Wang. 2023. [Agent instructs large language models to be general zero-shot reasoners](#).
- Gautier Dagan, Frank Keller, and Alex Lascarides. 2023. [Dynamic planning with a llm](#). *ArXiv preprint*, abs/2308.06391.
- Yuntian Deng, Kiran Prasad, Roland Fernandez, Paul Smolensky, Vishrav Chaudhary, and Stuart Shieber. 2023. [Implicit chain of thought reasoning via knowledge distillation](#).
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2023. [Chain-of-verification reduces hallucination in large language models](#). *ArXiv preprint*, abs/2309.11495.
- Shizhe Diao, Pengcheng Wang, Yong Lin, and Tong Zhang. 2023. [Active prompting with chain-of-thought for large language models](#). *ArXiv preprint*, abs/2302.12246.
- David Dohan, Winnie Xu, Aitor Lewkowycz, Jacob Austin, David Bieber, Raphael Gontijo Lopes, Yuhuai Wu, Henryk Michalewski, Rif A. Saurous, Jascha Sohl-Dickstein, Kevin Murphy, and Charles Sutton. 2022. [Language model cascades](#). *ArXiv preprint*, abs/2207.10342.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. 2023. [A survey for in-context learning](#). *ArXiv preprint*, abs/2301.00234.
- Qingxiu Dong, Ziwei Qin, Heming Xia, Tian Feng, Shoujie Tong, Haoran Meng, Lin Xu, Zhongyu Wei, Weidong Zhan, Baobao Chang, Sujian Li, Tianyu Liu, and Zhifang Sui. 2022. [Premise-based multimodal reasoning: Conditional inference on joint textual and visual clues](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 932–946, Dublin, Ireland. Association for Computational Linguistics.

840	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An image is worth 16x16 words: Transformers for image recognition at scale . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	898	Yao Fu, Hao-Chun Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. 2023b. Specializing smaller language models towards multi-step reasoning . In <i>International Conference on Machine Learning</i> .	899		900		901
844		902	Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. PAL: Program-aided language models . In <i>Proceedings of the 40th International Conference on Machine Learning</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 10764–10799. PMLR.	903		904		905
846		906		907		908		
847	Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multiagent debate . <i>ArXiv preprint</i> , abs/2305.14325.	909	Alfonso Emilio Gerevini. 2020. An introduction to the planning domain definition language (PDDL): book review . <i>Artif. Intell.</i> , 280:103221.	910		911		
848		912	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies . <i>Transactions of the Association for Computational Linguistics</i> , 9:346–361.	913		914		915
849	Dheeru Dua, Shivanshu Gupta, Sameer Singh, and Matt Gardner. 2022. Successive prompting for decomposing complex questions . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 1251–1265, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	916		917				
850		918	Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. 2023a. CRITIC: large language models can self-correct with tool-interactive critiquing . <i>ArXiv preprint</i> , abs/2305.11738.	919		920		921
851		922		923		924		925
852		926	Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu Chen. 2023b. Tora: A tool-integrated reasoning agent for mathematical problem solving .	927		928		929
853	Dheeru Dua, Sameer Singh, and Matt Gardner. 2020. Benefits of intermediate annotations in reading comprehension . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5627–5634, Online. Association for Computational Linguistics.	930	Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Linhao Yu, Yan Liu, Jiaxuan Li, Bojian Xiong, Deyi Xiong, et al. 2023. Evaluating large language models: A comprehensive survey. <i>arXiv preprint arXiv:2310.19736</i> .	931		932		933
854		934	Pranay Gupta and Manish Gupta. 2022. Newskvqa: Knowledge-aware news video question answering . In <i>Advances in Knowledge Discovery and Data Mining - 26th Pacific-Asia Conference, PAKDD 2022, Chengdu, China, May 16-19, 2022, Proceedings, Part III</i> , volume 13282 of <i>Lecture Notes in Computer Science</i> , pages 3–15. Springer.	935		936		937
855		938		939		940		941
856	Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.	942	Chengcheng Han, Xiaowei Du, Che Zhang, Yixin Lian, Xiang Li, Ming Gao, and Baoyuan Wang. 2023. Di-alCoT meets PPO: Decomposing and exploring reasoning paths in smaller language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 8055–8068, Singapore. Association for Computational Linguistics.	943		944		945
857		946		947		948		949
858		950	Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Luke Benson, Lucy Sun, Ekaterina Zubova, Yujie Qiao, Matthew Burtell, David Peng, Jonathan Fan, Yixin Liu, Brian Wong, Malcolm Sailor, Ansong Ni, Linyong Nan, Jungo Kasai, Tao Yu, Rui Zhang, Shafiq R. Joty, Alexander R. Fabbri, Wojciech Kryscinski, Xi Victoria Lin, Caiming	951		952		953
859								
860	Hao Fei, Bobo Li, Qian Liu, Lidong Bing, Fei Li, and Tat-Seng Chua. 2023. Reasoning implicit sentiment with chain-of-thought prompting . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 1171–1182. Association for Computational Linguistics.							
861								
862								
863								
864								
865								
866								
867								
868								
869								
870								
871								
872								
873								
874								
875								
876								
877								
878								
879								
880								
881								
882								
883								
884								
885								
886								
887								
888								
889								
890								
891								
892								
893								
894								
895								
896								
897								

954	Xiong, and Dragomir Radev. 2022. FOLIO: natural language reasoning with first-order logic . <i>ArXiv preprint</i> , abs/2209.00840.	2023 Conference on Empirical Methods in Natural Language Processing, pages 4902–4919, Singapore. Association for Computational Linguistics.	1010 1011 1012
957	Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen Wang, Daisy Wang, and Zhiting Hu. 2023a. Reasoning with language model is planning with world model . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 8154–8173, Singapore. Association for Computational Linguistics.	Cheng-Yu Hsieh, Si-An Chen, Chun-Liang Li, Yasuhisa Fujii, Alexander Ratner, Chen-Yu Lee, Ranjay Krishna, and Tomas Pfister. 2023a. Tool documentation enables zero-shot tool-usage with large language models . <i>ArXiv preprint</i> , abs/2308.00675.	1013 1014 1015 1016 1017
964	Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. 2023b. ToolkenGPT: Augmenting frozen language models with massive tools via tool embeddings . In <i>Thirty-seventh Conference on Neural Information Processing Systems, NeurIPS 2023</i> .	Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander J. Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023b. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes . <i>ArXiv preprint</i> , abs/2305.02301.	1018 1019 1020 1021 1022 1023
969	Hangfeng He, Hongming Zhang, and Dan Roth. 2023a. Rethinking with retrieval: Faithful large language model inference . <i>ArXiv preprint</i> , abs/2301.00303.	Hanxu Hu, Hongyuan Lu, Huajian Zhang, Wai Lam, and Yue Zhang. 2023a. Chain-of-symbol prompting elicits planning in large language models . <i>ArXiv preprint</i> , abs/2305.10276.	1024 1025 1026 1027
972	Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shuming Shi, and Xing Wang. 2023b. Exploring human-like translation strategy with large language models . <i>ArXiv preprint</i> , abs/2305.04118.	Mengkang Hu, Yao Mu, Xinmiao Yu, Mingyu Ding, Shiguang Wu, Wenqi Shao, Qiguang Chen, Bin Wang, Yu Qiao, and Ping Luo. 2023b. Tree-planner: Efficient close-loop task planning with large language models .	1028 1029 1030 1031 1032
977	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021a. Measuring massive multitask language understanding . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	Pengbo Hu, Ji Qi, Xingyu Li, Hong Li, Xinqi Wang, Bing Quan, Ruiyu Wang, and Yi Zhou. 2023c. Tree-of-mixed-thought: Combining fast and slow thinking for multi-hop visual reasoning . <i>ArXiv preprint</i> , abs/2308.09658.	1033 1034 1035 1036 1037
983	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. Measuring mathematical problem solving with the MATH dataset . In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .	Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2023a. Large language models can self-improve . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1051–1068, Singapore. Association for Computational Linguistics.	1038 1039 1040 1041 1042 1043
991	Namgyu Ho, Laura Schmid, and Se-Young Yun. 2023. Large language models are reasoning teachers . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 14852–14882. Association for Computational Linguistics.	Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2023b. Large language models cannot self-correct reasoning yet . <i>ArXiv preprint</i> , abs/2310.01798.	1044 1045 1046 1047 1048
998	Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. 2014. Learning to solve arithmetic word problems with verb categorization . In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 523–533, Doha, Qatar. Association for Computational Linguistics.	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023c. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions .	1049 1050 1051 1052 1053 1054
1005	Yifan Hou, Jiaoda Li, Yu Fei, Alessandro Stolfo, Wangchunshu Zhou, Guangtao Zeng, Antoine Bosselut, and Mrinmaya Sachan. 2023. Towards a mechanistic interpretation of multi-step reasoning capabilities of language models . In <i>Proceedings of the</i>	Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Cosmos QA: Machine reading comprehension with contextual commonsense reasoning . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2391–2401, Hong Kong, China. Association for Computational Linguistics.	1055 1056 1057 1058 1059 1060 1061 1062 1063

1064	Yue Huang, Jiawen Shi, Yuan Li, Chenrui Fan, Siyuan	Ben Mann, Sam McCandlish, Chris Olah, and Jared	1120
1065	Wu, Qihui Zhang, Yixin Liu, Pan Zhou, Yao Wan,	Kaplan. 2022. Language models (mostly) know what	1121
1066	Neil Zhenqiang Gong, and Lichao Sun. 2023d. Meta-	they know . <i>ArXiv preprint</i> , abs/2207.05221.	1122
1067	tool benchmark: Deciding whether to use tools and		
1068	which to use .		
1069	Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei	Ehud D. Karpas, Omri Abend, Yonatan Belinkov, Barak	1123
1070	Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu,	Lenz, Opher Lieber, Nir Ratner, Yoav Shoham, Hofit	1124
1071	Chuanheng Lv, Yikai Zhang, Jiayi Lei, Yao	Bata, Yoav Levine, Kevin Leyton-Brown, Dor Muhl-	1125
1072	Fu, Maosong Sun, and Junxian He. 2023e. C-	gay, Noam Rozen, Erez Schwartz, Gal Shachaf,	1126
1073	eval: A multi-level multi-discipline chinese eval-	Shai Shalev-Shwartz, Amnon Shashua, and Moshe	1127
1074	uation suite for foundation models . <i>ArXiv preprint</i> ,	Tenenholtz. 2022. Mrkl systems: A modular, neuro-	1128
1075	abs/2305.08322.	symbolic architecture that combines large language	1129
1076	Shima Imani, Liang Du, and Harsh Shrivastava. 2023.	models, external knowledge sources and discrete rea-	1130
1077	Mathprompter: Mathematical reasoning using large	soning . <i>ArXiv preprint</i> , abs/2205.00445.	1131
1078	language models . In <i>Proceedings of the The 61st An-</i>		
1079	<i>nual Meeting of the Association for Computational</i>	Uri Katz, Mor Geva, and Jonathan Berant. 2022. In-	1132
1080	<i>Linguistics: Industry Track, ACL 2023, Toronto,</i>	ferring implicit relations in complex questions with	1133
1081	<i>Canada, July 9-14, 2023</i> , pages 37–42. Association	language models . In <i>Findings of the Association</i>	1134
1082	for Computational Linguistics.	<i>for Computational Linguistics: EMNLP 2022</i> , pages	1135
1083	Raer Jack. 2023. Compression for agi . <i>Stanford MLSys</i> .	2548–2566, Abu Dhabi, United Arab Emirates. As-	1136
1084	Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko	sociation for Computational Linguistics.	1137
1085	Ishii, and Pascale Fung. 2023. Towards mitigat-		
1086	ing hallucination in large language models via self-	Muhammad Khalifa, Lajanugen Logeswaran, Moon-	1138
1087	reflection . <i>ArXiv preprint</i> , abs/2310.06271.	tae Lee, Honglak Lee, and Lu Wang. 2023.	1139
1088	Song Jiang, Zahra Shakeri, Aaron Chan, Maziar San-	Discriminator-guided multi-step reasoning with lan-	1140
1089	jabi, Hamed Firooz, Yinglong Xia, Bugra Akyildiz,	guage models . <i>ArXiv preprint</i> , abs/2305.14934.	1141
1090	Yizhou Sun, Jinchao Li, Qifan Wang, et al. 2023a.		
1091	Resprompt: Residual connection prompting advances	Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao	1142
1092	multi-step reasoning in large language models . <i>ArXiv</i>	Fu, Kyle Richardson, Peter Clark, and Ashish Sab-	1143
1093	<i>preprint</i> , abs/2310.04743.	harwal. 2023. Decomposed prompting: A modular	1144
1094	Weisen Jiang, Han Shi, Longhui Yu, Zhengying Liu,	approach for solving complex tasks . In <i>The Eleventh</i>	1145
1095	Yu Zhang, Zhenguo Li, and James T. Kwok. 2023b.	<i>International Conference on Learning Representa-</i>	1146
1096	Forward-backward reasoning in large language mod-	<i>tions, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> .	1147
1097	els for verification . <i>ArXiv preprint</i> , abs/2308.07758.	OpenReview.net.	1148
1098	Ziqi Jin and Wei Lu. 2023. Tab-cot: Zero-shot tabu-	Seungone Kim, Se Joo, Doyoung Kim, Joel Jang,	1149
1099	lar chain of thought . In <i>Findings of the Association</i>	Seonghyeon Ye, Jamin Shin, and Minjoon Seo.	1150
1100	<i>for Computational Linguistics: ACL 2023, Toronto,</i>	2023. The CoT collection: Improving zero-shot	1151
1101	<i>Canada, July 9-14, 2023</i> , pages 10259–10277. Asso-	and few-shot learning of language models via chain-	1152
1102	ciation for Computational Linguistics.	of-thought fine-tuning . In <i>Proceedings of the 2023</i>	1153
1103	Yan Junbing, Chengyu Wang, Taolin Zhang, Xiaofeng	<i>Conference on Empirical Methods in Natural Lan-</i>	1154
1104	He, Jun Huang, and Wei Zhang. 2023. From complex	<i>guage Processing</i> , pages 12685–12708, Singapore.	1155
1105	to simple: Unraveling the cognitive tree for reason-	Association for Computational Linguistics.	1156
1106	ing with small language models . In <i>Findings of the</i>		
1107	<i>Association for Computational Linguistics: EMNLP</i>	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-	1157
1108	2023, pages 12413–12425, Singapore. Association	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-	1158
1109	for Computational Linguistics.	guage models are zero-shot reasoners . In <i>NeurIPS</i> .	1159
1110	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom	Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish	1160
1111	Henighan, Dawn Drain, Ethan Perez, Nicholas	Sabharwal, Oren Etzioni, and Siena Dumas Ang.	1161
1112	Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli	2015. Parsing algebraic word problems into equa-	1162
1113	Tran-Johnson, Scott Johnston, Sheer El Showk, Andy	tions . <i>Transactions of the Association for Computa-</i>	1163
1114	Jones, Nelson Elhage, Tristan Hume, Anna Chen,	<i>tional Linguistics</i> , 3:585–597.	1164
1115	Yuntao Bai, Sam Bowman, Stanislav Fort, Deep		
1116	Ganguli, Danny Hernandez, Josh Jacobson, Jack-	Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate	1165
1117	son Kernion, Shauna Kravec, Liane Lovitt, Kam-	Kushman, and Hannaneh Hajishirzi. 2016. MAWPS:	1166
1118	al Ndousse, Catherine Olsson, Sam Ringer, Dario	A math word problem repository . In <i>Proceedings of</i>	1167
1119	Amodei, Tom Brown, Jack Clark, Nicholas Joseph,	<i>the 2016 Conference of the North American Chapter</i>	1168
		<i>of the Association for Computational Linguistics: Hu-</i>	1169
		<i>man Language Technologies</i> , pages 1152–1157, San	1170
		Diego, California. Association for Computational	1171
		Linguistics.	1172
		Yilun Kong, Jingqing Ruan, Yihong Chen, Bin Zhang,	1173
		Tianpeng Bao, Shiwei Shi, Guoqing Du, Xiaoru Hu,	1174
		Hangyu Mao, Ziyue Li, Xingyu Zeng, and Rui Zhao.	1175

1176	2023. Tptu-v2: Boosting task planning and tool usage of large language model-based agents in real-world systems.	Guy Gur-Ari, and Vedant Misra. 2022. Solving quantitative reasoning problems with language models. In <i>NeurIPS</i> .	1232
1177			1233
1178			1234
1179	Andrew Lampinen, Ishita Dasgupta, Stephanie Chan, Kory Mathewson, Mh Tessler, Antonia Creswell, James McClelland, Jane Wang, and Felix Hill. 2022. Can language models learn from explanations in context? In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 537–563, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	Jiangtong Li, Li Niu, and Liqing Zhang. 2022a. From representation to reasoning: Towards both evidence and commonsense reasoning for video question-answering. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 21241–21250. IEEE.	1235
1180			1236
1181			1237
1182			1238
1183			1239
1184			1240
1185			1241
1186			
1187	Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamile Lukosiute, Karina Nguyen, Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Robin Larson, Sam McCandlish, Sandipan Kundu, Saurav Kadavath, Shannon Yang, Thomas Henighan, Timothy Maxwell, Timothy Telleen-Lawton, Tristan Hume, Zac Hatfield-Dodds, Jared Kaplan, Jan Brauner, Samuel R. Bowman, and Ethan Perez. 2023. Measuring faithfulness in chain-of-thought reasoning. <i>ArXiv preprint</i> , abs/2307.13702.	Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. 2023a. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 19730–19742. PMLR.	1242
1188			1243
1189			1244
1190			1245
1191			1246
1192			1247
1193			1248
1194			1249
1195		Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai-Wei Chang, and Yejin Choi. 2023b. Symbolic chain-of-thought distillation: Small models can also "think" step-by-step. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 2665–2679. Association for Computational Linguistics.	1250
1196			1251
1197			1252
1198			1253
1199	Soochan Lee and Gunhee Kim. 2023. Recursion of thought: A divide-and-conquer approach to multi-context reasoning with language models. In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 623–658. Association for Computational Linguistics.		1254
1200			1255
1201			1256
1202			1257
1203		Minghao Li, Feifan Song, Bowen Yu, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023c. Api-bank: A benchmark for tool-augmented llms. <i>ArXiv preprint</i> , abs/2304.08244.	1258
1204			1259
1205	Bin Lei, Pei-Hung Lin, Chunhua Liao, and Caiwen Ding. 2023a. Boosting logical reasoning in large language models through a new framework: The graph of thought. <i>ArXiv preprint</i> , abs/2308.08614.		1260
1206			1261
1207		Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2022b. Contrastive decoding: Open-ended text generation as optimization. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	1262
1208			1263
1209	Deren Lei, Yaxi Li, Mingyu Wang, Vincent Yun, Emily Ching, Eslam Kamal, et al. 2023b. Chain of natural language inference for reducing large language model ungrounded hallucinations. <i>ArXiv preprint</i> , abs/2310.03951.		1264
1210			1265
1211			1266
1212			1267
1213		Xiaonan Li and Xipeng Qiu. 2023. Mot: Pre-thinking and recalling enable chatgpt to self-improve with memory-of-thoughts. <i>ArXiv preprint</i> , abs/2305.05181.	1268
1214	Jie Lei, Licheng Yu, Tamara Berg, and Mohit Bansal. 2020. What is more likely to happen next? video-and-language future event prediction. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 8769–8784, Online. Association for Computational Linguistics.		1269
1215			1270
1216			1271
1217		Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Lidong Bing, Shafiq R. Joty, and Soujanya Poria. 2023d. Chain of knowledge: A framework for grounding large language models with structured knowledge bases. <i>ArXiv preprint</i> , abs/2305.13269.	1272
1218			1273
1219			1274
1220			1275
1221	Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 19274–19286. PMLR.		1276
1222		Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, B. Chen, Jian-Guang Lou, and Weizhu Chen. 2022c. Making language models better reasoners with step-aware verifier. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	1277
1223			1278
1224			1279
1225			1280
1226			1281
1227			
1228	Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay V. Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur,	Yingcong Li, Kartik Sreenivasan, Angeliki Giannou, Dimitris S. Papailiopoulos, and Samet Oymak. 2023e. Dissecting chain-of-thought: A study on compositional in-context learning of mlps. <i>ArXiv preprint</i> , abs/2305.18869.	1282
1229			1283
1230			1284
1231			1285
			1286

1287	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris	with logical reasoning. In <i>Proceedings of the Twenty-</i>	1344
1288	Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian	<i>Ninth International Joint Conference on Artificial</i>	1345
1289	Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Ku-	<i>Intelligence, IJCAI 2020</i> , pages 3622–3628. ijcai.org.	1346
1290	mar, Benjamin Newman, Binhang Yuan, Bobby Yan,		
1291	Ce Zhang, Christian Cosgrove, Christopher D. Man-	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang,	1347
1292	ning, Christopher Ré, Diana Acosta-Navas, Drew A.	Hiroaki Hayashi, and Graham Neubig. 2023d. Pre-	1348
1293	Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak,	train, prompt, and predict: A systematic survey of	1349
1294	Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang,	prompting methods in natural language processing.	1350
1295	Keshav Santhanam, Laurel J. Orr, Lucia Zheng, Mert	<i>ACM Comput. Surv.</i> , 55(9):195:1–195:35.	1351
1296	Yüksekgönül, Mirac Suzgun, Nathan Kim, Neel		
1297	Guha, Niladri S. Chatterji, Omar Khattab, Peter	Tengxiao Liu, Qipeng Guo, Yuqing Yang, Xiangkun Hu,	1352
1298	Henderson, Qian Huang, Ryan Chi, Sang Michael	Yue Zhang, Xipeng Qiu, and Zheng Zhang. 2023e.	1353
1299	Xie, Shibani Santurkar, Surya Ganguli, Tatsunori	Plan, verify and switch: Integrated reasoning with	1354
1300	Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav	diverse X-of-thoughts. In <i>Proceedings of the 2023</i>	1355
1301	Chaudhary, William Wang, Xuechen Li, Yifan Mai,	<i>Conference on Empirical Methods in Natural Lan-</i>	1356
1302	Yuhui Zhang, and Yuta Koreeda. 2022. Holistic	<i>guage Processing</i> , pages 2807–2822, Singapore. As-	1357
1303	evaluation of language models. <i>ArXiv preprint,</i>	sociation for Computational Linguistics.	1358
1304	abs/2211.09110.		
1305	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang,	Tengxiao Liu, Qipeng Guo, Yuqing Yang, Xiangkun Hu,	1359
1306	Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and	Yue Zhang, Xipeng Qiu, and Zheng Zhang. 2023f.	1360
1307	Shuming Shi. 2023. Encouraging divergent thinking	Plan, verify and switch: Integrated reasoning with	1361
1308	in large language models through multi-agent debate.	diverse X-of-thoughts. In <i>Proceedings of the 2023</i>	1362
1309	<i>ArXiv preprint</i> , abs/2305.19118.	<i>Conference on Empirical Methods in Natural Lan-</i>	1363
		<i>guage Processing</i> , pages 2807–2822, Singapore. As-	1364
		sociation for Computational Linguistics.	1365
1310	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri	Jieyi Long. 2023. Large language model guided tree-of-	1366
1311	Edwards, Bowen Baker, Teddy Lee, Jan Leike, John	thought. <i>ArXiv preprint</i> , abs/2305.08291.	1367
1312	Schulman, Ilya Sutskever, and Karl Cobbe. 2023.		
1313	Let’s verify step by step.	Hongyuan Lu, Haoyang Huang, Dongdong Zhang, Hao-	1368
		ran Yang, Wai Lam, and Furu Wei. 2023a. Chain-	1369
1314	Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blun-	of-dictionary prompting elicits translation in large	1370
1315	som. 2017. Program induction by rationale genera-	language models. <i>ArXiv preprint</i> , abs/2305.06575.	1371
1316	tion: Learning to solve and explain algebraic word		
1317	problems. In <i>Proceedings of the 55th Annual Meet-</i>	Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-	1372
1318	<i>ing of the Association for Computational Linguistics</i>	Wei Chang, Song-Chun Zhu, Øyvind Taffjord, Peter	1373
1319	(<i>Volume 1: Long Papers</i>), pages 158–167, Vancouver,	Clark, and Ashwin Kalyan. 2022. Learn to explain:	1374
1320	Canada. Association for Computational Linguistics.	Multimodal reasoning via thought chains for science	1375
		question answering. In <i>NeurIPS</i> .	1376
1321	Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang,	Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu,	1377
1322	Mingu Lee, Roland Memisevic, and Hao Su. 2023.	Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark,	1378
1323	Deductive verification of chain-of-thought reasoning.	and Ashwin Kalyan. 2023b. Dynamic prompt learn-	1379
1324	In <i>Thirty-seventh Conference on Neural Information</i>	ing via policy gradient for semi-structured mathe-	1380
1325	<i>Processing Systems, NeurIPS 2023.</i>	matical reasoning. In <i>The Eleventh International</i>	1381
		<i>Conference on Learning Representations, ICLR 2023,</i>	1382
1326	Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi	<i>Kigali, Rwanda, May 1-5, 2023.</i> OpenReview.net.	1383
1327	Zhang, Joydeep Biswas, and Peter Stone. 2023a.		
1328	Llm+p: Empowering large language models with	Pan Lu, Liang Qiu, Wenhao Yu, Sean Welleck, and	1384
1329	optimal planning proficiency.	Kai-Wei Chang. 2023c. A survey of deep learn-	1385
		ing for mathematical reasoning. In <i>Proceedings</i>	1386
1330	Hanmeng Liu, Zhiyang Teng, Ruoxi Ning, Jian Liu,	<i>of the 61st Annual Meeting of the Association for</i>	1387
1331	Qiji Zhou, and Yue Zhang. 2023b. Glore: Evaluating	<i>Computational Linguistics (Volume 1: Long Papers),</i>	1388
1332	logical reasoning of large language models. <i>ArXiv</i>	<i>ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages	1389
1333	<i>preprint</i> , abs/2310.09107.	14605–14631. Association for Computational Lin-	1390
		guistics.	1391
1334	Jiacheng Liu, Ramakanth Pasunuru, Hannaneh Ha-	Yining Lu, Haoping Yu, and Daniel Khashabi. 2023d.	1392
1335	jishirzi, Yejin Choi, and Asli Celikyilmaz. 2023c.	Gear: Augmenting language models with generaliz-	1393
1336	Crystal: Introspective reasoners reinforced with self-	able and efficient tool resolution.	1394
1337	feedback. In <i>Proceedings of the 2023 Conference</i>		
1338	<i>on Empirical Methods in Natural Language Process-</i>	Man Luo, Shrinidhi Kumbhar, Ming shen, Mihir Parmar,	1395
1339	<i>ing</i> , pages 11557–11572, Singapore. Association for	Neeraj Varshney, Pratyay Banerjee, Somak Aditya,	1396
1340	Computational Linguistics.	and Chitta Baral. 2023. Towards logiglue: A brief	1397
1341	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang,	survey and a benchmark for analyzing logical reason-	1398
1342	Yile Wang, and Yue Zhang. 2020. Logiqa: A chal-	ing capabilities of language models.	1399
1343	lenge dataset for machine reading comprehension		

1400	Qianli Ma, Haotian Zhou, Tingkai Liu, Jianbo Yuan,	Dhabi, United Arab Emirates. Association for Com-	1457
1401	Pengfei Liu, Yang You, and Hongxia Yang. 2023.	putational Linguistics.	1458
1402	Let's reward step by step: Step-level reward model		
1403	as the navigators for reasoning.		
1404	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	Swaroop Mishra, Arindam Mitra, Neeraj Varshney,	1459
1405	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	Bhavdeep Sachdeva, Peter Clark, Chitta Baral, and	1460
1406	Nouha Dziri, Shrimai Prabhume, Yiming Yang,	Ashwin Kalyan. 2022b. NumGLUE: A suite of fun-	1461
1407	Shashank Gupta, Bodhisattwa Prasad Majumder,	damental yet challenging mathematical reasoning	1462
1408	Katherine Hermann, Sean Welleck, Amir Yazdan-	tasks. In <i>Proceedings of the 60th Annual Meeting of</i>	1463
1409	bakhsh, and Peter Clark. 2023. Self-refine: Itera-	<i>the Association for Computational Linguistics (Vol-</i>	1464
1410	tive refinement with self-feedback. In <i>Thirty-seventh</i>	<i>ume 1: Long Papers)</i> , pages 3505–3523, Dublin,	1465
1411	<i>Conference on Neural Information Processing Sys-</i>	Ireland. Association for Computational Linguistics.	1466
1412	<i>tems, NeurIPS 2023.</i>		
1413	Aman Madaan and Amir Yazdanbakhsh. 2022. Text	Shentong Mo and Miao Xin. 2023. Tree of uncertain	1467
1414	and patterns: For effective chain of thought, it takes	thoughts reasoning for large language models. <i>ArXiv</i>	1468
1415	two to tango. <i>ArXiv preprint</i> , abs/2209.07686.	<i>preprint</i> , abs/2309.07694.	1469
1416	Lucie Charlotte Magister, Jonathan Mallinson, Jakub	Ranjita Naik, Varun Chandrasekaran, Mert Yuksek-	1470
1417	Adámek, Eric Malmi, and Aliaksei Severyn. 2023.	gonul, Hamid Palangi, and Besmira Nushi. 2023.	1471
1418	Teaching small language models to reason. In <i>Pro-</i>	Diversity of thought improves reasoning abili-	1472
1419	<i>ceedings of the 61st Annual Meeting of the Asso-</i>	ties of large language models. <i>ArXiv preprint</i> ,	1473
1420	<i>ciation for Computational Linguistics (Volume 2:</i>	abs/2310.07088.	1474
1421	<i>Short Papers)</i> , ACL 2023, Toronto, Canada, July		
1422	9-14, 2023, pages 1773–1781. Association for Com-	Deepak Nathani, David Wang, Liangming Pan, and	1475
1423	putational Linguistics.	William Wang. 2023. MAF: Multi-aspect feedback	1476
1424	Ana Marasovic, Iz Beltagy, Doug Downey, and Matthew	for improving reasoning in large language models.	1477
1425	Peters. 2022. Few-shot self-rationalization with nat-	In <i>Proceedings of the 2023 Conference on Empiri-</i>	1478
1426	ural language prompts. In <i>Findings of the Associa-</i>	<i>cal Methods in Natural Language Processing</i> , pages	1479
1427	<i>tion for Computational Linguistics: NAACL 2022</i> ,	6591–6616, Singapore. Association for Computa-	1480
1428	pages 410–424, Seattle, United States. Association	tional Linguistics.	1481
1429	for Computational Linguistics.		
1430	William Merrill and Ashish Sabharwal. 2023. The	Xuefei Ning, Zinan Lin, Zixuan Zhou, Huazhong Yang,	1482
1431	expressive power of transformers with chain of	and Yu Wang. 2023. Skeleton-of-thought: Large	1483
1432	thought.	language models can do parallel decoding. <i>ArXiv</i>	1484
1433	Ning Miao, Yee Whye Teh, and Tom Rainforth.	<i>preprint</i> , abs/2307.15337.	1485
1434	2023. Selfcheck: Using llms to zero-shot check	Sean O'Brien and Mike Lewis. 2023. Contrastive de-	1486
1435	their own step-by-step reasoning. <i>ArXiv preprint</i> ,	coding improves reasoning in large language models.	1487
1436	abs/2308.00436.	<i>ArXiv preprint</i> , abs/2309.09117.	1488
1437	Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su.	Theo Olausson, Alex Gu, Ben Lipkin, Cedegao Zhang,	1489
1438	2020. A diverse corpus for evaluating and developing	Armando Solar-Lezama, Joshua Tenenbaum, and	1490
1439	English math word problem solvers. In <i>Proceedings</i>	Roger Levy. 2023. LINC: A neurosymbolic approach	1491
1440	<i>of the 58th Annual Meeting of the Association for</i>	for logical reasoning by combining language models	1492
1441	<i>Computational Linguistics</i> , pages 975–984, Online.	with first-order logic provers. In <i>Proceedings of the</i>	1493
1442	Association for Computational Linguistics.	<i>2023 Conference on Empirical Methods in Natural</i>	1494
1443	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish	<i>Language Processing</i> , pages 5153–5176, Singapore.	1495
1444	Sabharwal. 2018. Can a suit of armor conduct elec-	Association for Computational Linguistics.	1496
1445	tricity? a new dataset for open book question an-	OpenAI. 2023. GPT-4 technical report. <i>ArXiv preprint</i> ,	1497
1446	swering. In <i>Proceedings of the 2018 Conference on</i>	abs/2303.08774.	1498
1447	<i>Empirical Methods in Natural Language Processing</i> ,	Liangming Pan, Alon Albalak, Xinyi Wang, and	1499
1448	pages 2381–2391, Brussels, Belgium. Association	William Wang. 2023. Logic-LM: Empowering large	1500
1449	for Computational Linguistics.	language models with symbolic solvers for faithful	1501
1450	Swaroop Mishra, Matthew Finlayson, Pan Lu, Leonard	logical reasoning. In <i>Findings of the Association</i>	1502
1451	Tang, Sean Welleck, Chitta Baral, Tanmay Rajpuro-	<i>for Computational Linguistics: EMNLP 2023</i> , pages	1503
1452	hit, Oyvind Taffjord, Ashish Sabharwal, Peter Clark,	3806–3824, Singapore. Association for Computa-	1504
1453	and Ashwin Kalyan. 2022a. LILA: A unified bench-	tional Linguistics.	1505
1454	mark for mathematical reasoning. In <i>Proceedings of</i>	Bhargavi Paranjape, Scott Lundberg, Sameer Singh,	1506
1455	<i>the 2022 Conference on Empirical Methods in Nat-</i>	Hannaneh Hajishirzi, Luke Zettlemoyer, and	1507
1456	<i>ural Language Processing</i> , pages 5807–5832, Abu	Marco Tulio Ribeiro. 2023. Art: Automatic multi-	1508
		step reasoning and tool-use for large language mod-	1509
		els.	1510

1511	Aaron Parisi, Yao Zhao, and Noah Fiedel. 2022. Talm: Tool augmented language models . <i>ArXiv preprint</i> , abs/2205.12255.	1568
1512		1569
1513		
1514	Jae Sung Park, Chandra Bhagavatula, Roozbeh Motaghi, Ali Farhadi, and Yejin Choi. 2020. Visual-comet: Reasoning about the dynamic context of a still image . In <i>Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part V</i> , volume 12350 of <i>Lecture Notes in Computer Science</i> , pages 508–524. Springer.	1570
1515		1571
1516		1572
1517		1573
1518		1574
1519		1575
1520		1576
1521	Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP models really able to solve simple math word problems? In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2080–2094, Online. Association for Computational Linguistics.	1577
1522		1578
1523		1579
1524		1580
1525		
1526		1581
1527		1582
1528	Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. 2023. REFINER: reasoning feedback on intermediate representations . <i>ArXiv preprint</i> , abs/2304.01904.	1583
1529		1584
1530		1585
1531		1586
1532		1587
1533	Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding multimodal large language models to the world . <i>ArXiv preprint</i> , abs/2306.14824.	1588
1534		1589
1535		1590
1536		1591
1537		
1538	Silviu Pitis, Michael R. Zhang, Andrew Wang, and Jimmy Ba. 2023. Boosted prompt ensembles for large language models . <i>ArXiv preprint</i> , abs/2304.05970.	1592
1539		1593
1540		1594
1541		1595
1542		1596
1543		1597
1544		1598
1545	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah Smith, and Mike Lewis. 2023. Measuring and narrowing the compositionality gap in language models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 5687–5711, Singapore. Association for Computational Linguistics.	1599
1546		1600
1547		1601
1548	Jingyuan Qi, Zhiyang Xu, Ying Shen, Minqian Liu, Di Jin, Qifan Wang, and Lifu Huang. 2023. The art of SOCRATIC QUESTIONING: Recursive thinking with large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 4177–4199, Singapore. Association for Computational Linguistics.	1602
1549		1603
1550		1604
1551		1605
1552		
1553		1606
1554		1607
1555	Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang, and Huajun Chen. 2023. Reasoning with language model prompting: A survey . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , ACL 2023, Toronto, Canada, July 9-14, 2023, pages 5368–5393. Association for Computational Linguistics.	1608
1556		1609
1557		1610
1558		1611
1559		1612
1560		
1561		1613
1562		1614
1563	Libo Qin, Qiguang Chen, Fuxuan Wei, Shijue Huang, and Wanxiang Che. 2023a. Cross-lingual prompting: Improving zero-shot chain-of-thought reasoning across languages . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language</i>	1615
1564		1616
1565		1617
1566		
1567		
	<i>Processing</i> , pages 2695–2709, Singapore. Association for Computational Linguistics.	
	Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023b. Toolllm: Facilitating large language models to master 16000+ real-world apis . <i>ArXiv preprint</i> , abs/2307.16789.	
	Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. 2020. Pre-trained models for natural language processing: A survey . <i>ArXiv preprint</i> , abs/2003.08271.	
	Ansh Radhakrishnan, Karina Nguyen, Anna Chen, Carol Chen, Carson Denison, Danny Hernandez, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamile Lukosiute, Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Sam McCandlish, Sheer El Showk, Tamera Lanham, Tim Maxwell, Venkatesa Chandrasekaran, Zac Hatfield-Dodds, Jared Kaplan, Jan Brauner, Samuel R. Bowman, and Ethan Perez. 2023. Question decomposition improves the faithfulness of model-generated reasoning . <i>ArXiv preprint</i> , abs/2307.11768.	
	Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019a. Explain yourself! leveraging language models for commonsense reasoning . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4932–4942, Florence, Italy. Association for Computational Linguistics.	
	Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019b. Explain yourself! leveraging language models for commonsense reasoning . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4932–4942, Florence, Italy. Association for Computational Linguistics.	
	Hannah Rashkin, Maarten Sap, Emily Allaway, Noah A. Smith, and Yejin Choi. 2018. Event2Mind: Commonsense inference on events, intents, and reactions . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 463–473, Melbourne, Australia. Association for Computational Linguistics.	
	Subhro Roy and Dan Roth. 2015. Solving general arithmetic word problems . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 1743–1752, Lisbon, Portugal. Association for Computational Linguistics.	
	Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu, Tianpeng Bao, Guoqing Du, Shiwei Shi, Hangyu Mao, Xingyu Zeng, and Rui Zhao. 2023. Tptu: Task planning and tool usage of large language model-based ai agents .	

1623	Abulhair Saparov and He He. 2023. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	1681
1624		1682
1625		
1626		
1627		
1628	Teven Le Scao, Angela Fan, Christopher Akiki, El-	
1629	lie Pavlick, Suzana Ilic, Daniel Hesslow, Roman	
1630	Castagné, Alexandra Sasha Luccioni, François Yvon,	
1631	Matthias Gallé, Jonathan Tow, Alexander M. Rush,	
1632	Stella Biderman, Albert Webson, Pawan Sasanka Am-	
1633	manamanchi, Thomas Wang, Benoît Sagot, Niklas	
1634	Muennighoff, Albert Villanova del Moral, Olatunji	
1635	Ruwase, Rachel Bawden, Stas Bekman, Angelina	
1636	McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile	
1637	Saulnier, Samson Tan, Pedro Ortiz Suarez, Vic-	
1638	tor Sanh, Hugo Laurençon, Yacine Jernite, Julien	
1639	Launay, Margaret Mitchell, Colin Raffel, Aaron	
1640	Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri	
1641	Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg	
1642	Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue,	
1643	Christopher Klam, Colin Leong, Daniel van Strien,	
1644	David Ifeoluwa Adelani, and et al. 2022. BLOOM: A 176b-parameter open-access multilingual language model . <i>ArXiv preprint</i> , abs/2211.05100.	
1645		
1646		
1647	Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo.	
1648	2023. Are emergent abilities of large language models a mirage? In <i>Thirty-seventh Conference on Neural Information Processing Systems, NeurIPS 2023</i> .	
1649		
1650		
1651	Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta	
1652	Raileanu, Maria Lomeli, Eric Hambro, Luke Zettle-	
1653	moyer, Nicola Cancedda, and Thomas Scialom. 2023.	
1654	Toolformer: Language models can teach themselves to use tools . In <i>Thirty-seventh Conference on Neural Information Processing Systems, NeurIPS 2023</i> .	
1655		
1656		
1657	Bilgehan Sel, Ahmad Al-Tawaha, Vanshaj Khattar,	
1658	Lu Wang, Ruoxi Jia, and Ming Jin. 2023. Algorithm of thoughts: Enhancing exploration of ideas in large language models . <i>ArXiv preprint</i> , abs/2308.10379.	
1659		
1660		
1661	Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie	
1662	Huang, Nan Duan, and Weizhu Chen. 2023a. Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9248–9274, Singapore. Association for Computational Linguistics.	
1663		
1664		
1665		
1666		
1667		
1668	Zhihong Shao, Yeyun Gong, Yelong Shen, Min-	
1669	lie Huang, Nan Duan, and Weizhu Chen. 2023b. Synthetic prompting: Generating chain-of-thought demonstrations for large language models . <i>ArXiv preprint</i> , abs/2302.00618.	
1670		
1671		
1672		
1673	Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li,	
1674	Weiming Lu, and Yueting Zhuang. 2023a. Hugging-GPT: Solving AI tasks with chatGPT and its friends in hugging face . In <i>Thirty-seventh Conference on Neural Information Processing Systems, NeurIPS 2023</i> .	
1675		
1676		
1677		
1678		
1679	Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang,	
1680	Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng Li,	
	and Yueting Zhuang. 2023b. Taskbench: Benchmarking large language models for task automation .	1681
		1682
	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang,	1683
	Suraj Srivats, Soroush Vosoughi, Hyung Won Chung,	1684
	Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das,	1685
	and Jason Wei. 2023. Language models are multi-lingual chain-of-thought reasoners . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	1686
		1687
		1688
		1689
		1690
	Noah Shinn, Federico Cassano, Ashwin Gopinath,	1691
	Karthik R Narasimhan, and Shunyu Yao. 2023. Re-reflexion: language agents with verbal reinforcement learning . In <i>Thirty-seventh Conference on Neural Information Processing Systems, NeurIPS 2023</i> .	1692
		1693
		1694
		1695
	Kumar Shridhar, Harsh Jhamtani, Hao Fang, Benjamin	1696
	Van Durme, Jason Eisner, and Patrick Xia. 2023. Screws: A modular framework for reasoning with revisions . <i>ArXiv preprint</i> , abs/2309.13075.	1697
		1698
		1699
	Kashun Shum, Shizhe Diao, and Tong Zhang. 2023. Automatic prompt augmentation and selection with chain-of-thought from labeled data . <i>ArXiv preprint</i> , abs/2302.12822.	1700
		1701
		1702
		1703
	Aaro,hi Srivastava, Abhinav Rastogi, Abhishek Rao,	1704
	Abu Awal Md Shoeb, Abubakar Abid, Adam	1705
	Fisch, Adam R. Brown, Adam Santoro, Aditya	1706
	Gupta, Adrià Garriga-Alonso, Agnieszka Kluska,	1707
	Aitor Lewkowycz, Akshat Agarwal, Alethea Power,	1708
	Alex Ray, Alex Warstadt, Alexander W. Kocurek,	1709
	Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Par-	1710
	rish, Allen Nie, Aman Hussain, Amanda Askill,	1711
	Amanda Dsouza, Ameet Rahane, Anantharaman S.	1712
	Iyer, Anders Andreassen, Andrea Santilli, Andreas	1713
	Stuhlmüller, Andrew M. Dai, Andrew La, Andrew K.	1714
	Lampinen, Andy Zou, Angela Jiang, Angelica Chen,	1715
	Anh Vuong, Animesh Gupta, Anna Gottardi, Anto-	1716
	nio Norelli, Anu Venkatesh, Arash Gholamidavoodi,	1717
	Arfa Tabassum, Arul Menezes, Arun Kirubarajan,	1718
	Asher Mullokandov, Ashish Sabharwal, Austin Her-	1719
	rick, Avia Efrat, Aykut Erdem, Ayla Karakas, and	1720
	et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models . <i>ArXiv preprint</i> , abs/2206.04615.	1721
		1722
		1723
	Haotian Sun, Yuchen Zhuang, Ling kai Kong, Bo Dai,	1724
	and Chao Zhang. 2023. Adaplaner: Adaptive planning from feedback with language models . In <i>Thirty-seventh Conference on Neural Information Processing Systems, NeurIPS 2023</i> .	1725
		1726
		1727
		1728
	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Se-	1729
	bastian Gehrmann, Yi Tay, Hyung Won Chung,	1730
	Aakanksha Chowdhery, Quoc V. Le, Ed Chi, Denny	1731
	Zhou, and Jason Wei. 2023. Challenging big-bench tasks and whether chain-of-thought can solve them . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 13003–13051. Association for Computational Linguistics.	1732
		1733
		1734
		1735
		1736
		1737

- Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. 2021. [ProofWriter: Generating implications, proofs, and abductive statements over natural language](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3621–3634, Online. Association for Computational Linguistics.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A question answering challenge targeting commonsense knowledge](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alon Talmor, Ori Yoran, Ronan Le Bras, Chandra Bhagavatula, Yoav Goldberg, Yejin Choi, and Jonathan Berant. 2021. [Commonsenseqa 2.0: Exposing the limits of AI through gamification](#). In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*.
- Xiaojuan Tang, Zilong Zheng, Jiaqi Li, Fanxu Meng, Song-Chun Zhu, Yitao Liang, and Muhan Zhang. 2023. [Large language models are in-context semantic reasoners rather than symbolic reasoners](#). *ArXiv preprint*, abs/2305.14825.
- Qingyuan Tian, Hanlun Zhu, Lei Wang, Yang Li, and Yunshi Lan. 2023. [R³ prompting: Review, rephrase and resolve for chain-of-thought reasoning in large language models under noisy context](#).
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. [Llama: Open and efficient foundation language models](#). *ArXiv preprint*, abs/2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas
- Scialom. 2023b. [Llama 2: Open foundation and fine-tuned chat models](#). *ArXiv preprint*, abs/2307.09288.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. [Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2023, Toronto, Canada, July 9-14, 2023, pages 10014–10037. Association for Computational Linguistics.
- Rasul Tutunov, Antoine Grosnit, Juliusz Ziomek, Jun Wang, and Haitham Bou-Ammar. 2023. [Why can large language models generate correct chain-of-thoughts?](#) *ArXiv preprint*, abs/2310.13571.
- Jonathan Uesato, Nate Kushman, Ramana Kumar, H. Francis Song, Noah Y. Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. 2022. [Solving math word problems with process- and outcome-based feedback](#). *ArXiv preprint*, abs/2211.14275.
- Xingchen Wan, Ruoxi Sun, Hanjun Dai, Serkan Ö. Arik, and Tomas Pfister. 2023. [Better zero-shot reasoning with self-adaptive prompting](#). In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 3493–3514. Association for Computational Linguistics.
- Boshi Wang, Xiang Deng, and Huan Sun. 2022. [Iteratively prompt pre-trained language models for chain of thought](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2714–2730, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. 2023a. [Towards understanding chain-of-thought prompting: An empirical study of what matters](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2023, Toronto, Canada, July 9-14, 2023, pages 2717–2739. Association for Computational Linguistics.
- Cunxiang Wang, Shuailong Liang, Yue Zhang, Xiaonan Li, and Tian Gao. 2019. [Does it make sense? and why? a pilot study for sense making and explanation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4020–4026, Florence, Italy. Association for Computational Linguistics.
- Haotian Wang, Xiyuan Du, Weijiang Yu, Qianglong Chen, Kun Zhu, Zheng Chu, Lian Yan, and Yi Guan. 2023b. [Apollo’s oracle: Retrieval-augmented reasoning in multi-agent debates](#).
- Jianing Wang, Qiushi Sun, Nuo Chen, Xiang Li, and Ming Gao. 2023c. [Boosting language models reasoning with chain-of-knowledge prompting](#). *ArXiv preprint*, abs/2306.06427.

- Jinyuan Wang, Junlong Li, and Hai Zhao. 2023d. [Self-prompted chain-of-thought on large language models for open-domain multi-hop reasoning](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2717–2731, Singapore. Association for Computational Linguistics. 1910
- Keheng Wang, Feiyu Duan, Sirui Wang, Peiguang Li, Yunsen Xian, Chuantao Yin, Wenge Rong, and Zhang Xiong. 2023e. [Knowledge-driven cot: Exploring faithful reasoning in llms for knowledge-intensive question answering](#). 1911
- Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun. 2023f. [Label words are anchors: An information flow perspective for understanding in-context learning](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9840–9855, Singapore. Association for Computational Linguistics. 1912
- Lei Wang, Yi Hu, Jiabang He, Xing Xu, Ning Liu, Hui Liu, and Heng Tao Shen. 2023g. [T-sciq: Teaching multimodal chain-of-thought reasoning via large language model signals for science question answering](#). *ArXiv preprint*, abs/2305.03453. 1913
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji-Rong Wen. 2023h. [A survey on large language model based autonomous agents](#). *ArXiv preprint*, abs/2308.11432. 1914
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023i. [Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2023, Toronto, Canada, July 9-14, 2023, pages 2609–2634. Association for Computational Linguistics. 1915
- Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao, Bing Yin, and Xiang Ren. 2023j. [Scott: Self-consistent chain-of-thought distillation](#). In *Annual Meeting of the Association for Computational Linguistics*. 1916
- Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2023k. [Mint: Evaluating llms in multi-turn interaction with tools and language feedback](#). 1917
- Xinyi Wang, Lucas Caccia, Oleksiy Ostapenko, Xingdi Yuan, and Alessandro Sordoni. 2023l. [Guiding language model reasoning with planning tokens](#). 1918
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023m. [Self-consistency improves chain of thought reasoning in language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. 1919
- Yiming Wang, Zhuosheng Zhang, and Rui Wang. 2023n. [Element-aware summarization with large language models: Expert-aligned evaluation and chain-of-thought method](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2023, Toronto, Canada, July 9-14, 2023, pages 8640–8665. Association for Computational Linguistics. 1920
- Yuqing Wang and Yun Zhao. 2023. [TRAM: benchmarking temporal reasoning for large language models](#). *ArXiv preprint*, abs/2310.00835. 1921
- Zhaoyang Wang, Shaohan Huang, Yuxuan Liu, Jiahai Wang, Minghui Song, Zihan Zhang, Haizhen Huang, Furu Wei, Weiwei Deng, Feng Sun, and Qi Zhang. 2023o. [Democratizing reasoning ability: Tailored learning from large language model](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1948–1966, Singapore. Association for Computational Linguistics. 1922
- 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. [Emergent abilities of large language models](#). *Trans. Mach. Learn. Res.*, 2022. 1923
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022b. [Chain-of-thought prompting elicits reasoning in large language models](#). In *NeurIPS*. 1924
- Yixuan Weng, Minjun Zhu, Shizhu He, Kang Liu, and Jun Zhao. 2022. [Large language models are reasoners with self-verification](#). *ArXiv preprint*, abs/2212.09561. 1925
- Bo Wu, Shoubin Yu, Zhenfang Chen, Josh Tenenbaum, and Chuang Gan. 2021. [STAR: A benchmark for situated reasoning in real-world videos](#). In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*. 1926
- Haoyi Wu, Wenyang Hui, Yezeng Chen, Weiqi Wu, Kewei Tu, and Yi Zhou. 2023a. [Conic10K: A challenging math problem understanding and reasoning dataset](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6444–6458, Singapore. Association for Computational Linguistics. 1927
- Skyler Wu, Eric Meng Shen, Charumathi Badrinath, Jiaqi Ma, and Himabindu Lakkaraju. 2023b. [Analyzing chain-of-thought prompting in large language models via gradient-based feature attributions](#). *ArXiv preprint*, abs/2307.13339. 1928

1966	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huan, and Tao Gui. 2023. The rise and potential of large language model based agents: A survey . <i>ArXiv preprint</i> , abs/2309.07864.	<i>Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	2022 2023
1967			
1968			
1969		Yao Yao, Zuchao Li, and Hai Zhao. 2023d. Beyond chain-of-thought, effective graph-of-thought reasoning in large language models . <i>ArXiv preprint</i> , abs/2305.16582.	2024 2025 2026 2027
1970			
1971			
1972			
1973			
1974			
1975		Xi Ye, Qiaochu Chen, Isil Dillig, and Greg Durrett. 2023a. SatLM: Satisfiability-aided language models using declarative prompting . In <i>Thirty-seventh Conference on Neural Information Processing Systems, NeurIPS 2023</i> .	2028 2029 2030 2031 2032
1976			
1977	Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. Next-qa: Next phase of question-answering to explaining temporal actions . In <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021</i> , pages 9777–9786. Computer Vision Foundation / IEEE.		
1978			
1979			
1980			
1981		Xi Ye and Greg Durrett. 2022. The unreliability of explanations in few-shot in-context learning . <i>ArXiv preprint</i> , abs/2205.03401.	2033 2034 2035
1982			
1983	Weijia Xu, Andrzej Banburski-Fahey, and Nebojsa Jojic. 2023. Reprompting: Automated chain-of-thought prompt inference through gibbs sampling .	Xi Ye and Greg Durrett. 2023. Explanation selection using unlabeled data for in-context learning . <i>ArXiv preprint</i> , abs/2302.04813.	2036 2037 2038
1984			
1985			
1986	Tianci Xue, Ziqi Wang, Zhenhailong Wang, Chi Han, Pengfei Yu, and Heng Ji. 2023. RCOT: detecting and rectifying factual inconsistency in reasoning by reversing chain-of-thought . <i>ArXiv preprint</i> , abs/2305.11499.	Yunhu Ye, Binyuan Hui, Min Yang, Binhua Li, Fei Huang, and Yongbin Li. 2023b. Large language models are versatile decomposers: Decomposing evidence and questions for table-based reasoning . In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2023, Taipei, Taiwan, July 23-27, 2023</i> , pages 174–184. ACM.	2039 2040 2041 2042 2043 2044 2045 2046
1987			
1988			
1989			
1990			
1991	Hui Yang, Sifu Yue, and Yunzhong He. 2023a. Auto-gpt for online decision making: Benchmarks and additional opinions . <i>ArXiv preprint</i> , abs/2306.02224.		
1992			
1993			
1994	Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. 2023b. MM-REACT: prompting chatgpt for multimodal reasoning and action . <i>ArXiv preprint</i> , abs/2303.11381.	Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B. Tenenbaum. 2020. CLEVRER: collision events for video representation and reasoning . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	2047 2048 2049 2050 2051 2052 2053
1995			
1996			
1997			
1998			
1999	Zonglin Yang, Li Dong, Xinya Du, Hao Cheng, Erik Cambria, Xiaodong Liu, Jianfeng Gao, and Furu Wei. 2022. Language models as inductive reasoners . <i>ArXiv preprint</i> , abs/2212.10923.	Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023. Do large language models know what they don't know? In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 8653–8665. Association for Computational Linguistics.	2054 2055 2056 2057 2058 2059 2060
2000			
2001			
2002			
2003	Zonglin Yang, Xinya Du, Rui Mao, Jinjie Ni, and Erik Cambria. 2023c. Logical reasoning over natural language as knowledge representation: A survey . <i>ArXiv preprint</i> , abs/2303.12023.		
2004			
2005			
2006			
2007	Fanglong Yao, Changyuan Tian, Jintao Liu, Zequn Zhang, Qing Liu, Li Jin, Shuchao Li, Xiaoyu Li, and Xian Sun. 2023a. Thinking like an expert:multimodal hypergraph-of-thought (hot) reasoning to boost foundation modals .	Ori Yoran, Tomer Wolfson, Ben Bogin, Uri Katz, Daniel Deutch, and Jonathan Berant. 2023. Answering questions by meta-reasoning over multiple chains of thought . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 5942–5966, Singapore. Association for Computational Linguistics.	2061 2062 2063 2064 2065 2066 2067
2008			
2009			
2010			
2011			
2012	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R Narasimhan. 2023b. Tree of thoughts: Deliberate problem solving with large language models . In <i>Thirty-seventh Conference on Neural Information Processing Systems, NeurIPS 2023</i> .	Fei Yu, Hongbo Zhang, and Benyou Wang. 2023a. Nature language reasoning, A survey . <i>ArXiv preprint</i> , abs/2303.14725.	2068 2069 2070
2013			
2014			
2015			
2016			
2017			
2018	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023c. React: Synergizing reasoning and acting in language models . In <i>The Eleventh International</i>	Junchi Yu, Ran He, and Rex Ying. 2023b. Thought propagation: An analogical approach to complex reasoning with large language models . <i>ArXiv preprint</i> , abs/2310.03965.	2071 2072 2073 2074
2019			
2020			
2021			

2075	Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng.	Yifan Zhang, Jingqin Yang, Yang Yuan, and An-	2130
2076	2020. Reclor: A reading comprehension dataset re-	drew Chi-Chih Yao. 2023e. Cumulative reason-	2131
2077	quiring logical reasoning . In <i>8th International Con-</i>	ing with large language models . <i>ArXiv preprint</i> ,	2132
2078	<i>ference on Learning Representations, ICLR 2020</i> ,	abs/2308.04371 .	2133
2079	<i>Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenRe-		
2080	view.net .		
2081	Xiao Yu, Baolin Peng, Michel Galley, Jianfeng Gao,	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,	2134
2082	and Zhou Yu. 2023c. Teaching language models to	Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,	2135
2083	self-improve through interactive demonstrations .	Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei	2136
2084		Bi, Freda Shi, and Shuming Shi. 2023f. Siren’s song	2137
2085	Zihan Yu, Liang He, Zhen Wu, Xinyu Dai, and Jia-	in the AI ocean: A survey on hallucination in large	2138
2086	jun Chen. 2023d. Towards better chain-of-thought	language models . <i>ArXiv preprint</i> , abs/2309.01219 .	2139
2087	prompting strategies: A survey .		
2088	Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D.	Zhebin Zhang, Xinyu Zhang, Yuanhang Ren, Saijiang	2140
2089	Goodman. 2022. Star: Bootstrapping reasoning with	Shi, Meng Han, Yongkang Wu, Ruofei Lai, and Zhao	2141
	reasoning . In <i>NeurIPS</i> .	Cao. 2023g. IAG: Induction-augmented generation	2142
2090	Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin	framework for answering reasoning questions . In	2143
2091	Choi. 2019. From recognition to cognition: Visual	<i>Proceedings of the 2023 Conference on Empirical</i>	2144
2092	commonsense reasoning . In <i>IEEE Conference on</i>	<i>Methods in Natural Language Processing</i> , pages 1–	2145
2093	<i>Computer Vision and Pattern Recognition, CVPR</i>	14, Singapore. Association for Computational Lin-	2146
2094	<i>2019, Long Beach, CA, USA, June 16-20, 2019</i> , pages	<i>guistics</i> .	2147
2095	6720–6731. Computer Vision Foundation / IEEE.		
2096	Bowen Zhang, Kehua Chang, and Chunping Li. 2023a.	Zhuosheng Zhang and Aston Zhang. 2023. You only	2148
2097	Cot-bert: Enhancing unsupervised sentence repre-	look at screens: Multimodal chain-of-action agents .	2149
2098	sentation through chain-of-thought . <i>ArXiv preprint</i> ,		
2099	abs/2309.11143 .		
2100	Hugh Zhang and David C. Parkes. 2023. Chain-of-	Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex	2150
2101	thought reasoning is a policy improvement operator .	Smola. 2023h. Automatic chain of thought prompt-	2151
2102	Jun Zhang, Jue Wang, Huan Li, Lidan Shou, Ke Chen,	ing in large language models . In <i>The Eleventh In-</i>	2152
2103	Gang Chen, and Sharad Mehrotra. 2023b. Draft	<i>ternational Conference on Learning Representations</i> ,	2153
2104	& verify: Lossless large language model accelera-	<i>ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . Open-	2154
2105	tion via self-speculative decoding . <i>ArXiv preprint</i> ,	Review.net .	2155
2106	abs/2309.08168 .		
2107	Muru Zhang, Ofir Press, William Merrill, Alisa	Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao,	2156
2108	Liu, and Noah A. Smith. 2023c. How language	George Karypis, and Alex Smola. 2023i. Multi-	2157
2109	model hallucinations can snowball . <i>ArXiv preprint</i> ,	modal chain-of-thought reasoning in language mod-	2158
2110	abs/2305.13534 .	els . <i>ArXiv preprint</i> , abs/2302.00923 .	2159
2111	Sarah J. Zhang, Reece Shuttleworth, Derek Austin,	Ruochen Zhao, Xingxuan Li, Shafiq Joty, Chengwei	2160
2112	Yann Hicke, Leonard Tang, Sathwik Karnik, Dar-	Qin, and Lidong Bing. 2023a. Verify-and-edit: A	2161
2113	nell Granberry, and Iddo Drori. 2022a. A dataset and	knowledge-enhanced chain-of-thought framework .	2162
2114	benchmark for automatically answering and gener-	In <i>Proceedings of the 61st Annual Meeting of the</i>	2163
2115	ating machine learning final exams . <i>ArXiv preprint</i> ,	<i>Association for Computational Linguistics (Volume</i>	2164
2116	abs/2206.05442 .	<i>1: Long Papers)</i> , <i>ACL 2023, Toronto, Canada, July</i>	2165
2117	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	9-14, 2023, pages 5823–5840. Association for Com-	2166
2118	Artetxe, Moya Chen, Shuohui Chen, Christopher	<i>putational Linguistics</i> .	2167
2119	Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin,	Wayne Xin Zhao, Kun Zhou, Zheng Gong, Beichen	2168
2120	Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shus-	Zhang, Yuanhang Zhou, Jing Sha, Zhigang Chen,	2169
2121	ter, Daniel Simig, Punit Singh Koura, Anjali Sridhar,	Shijin Wang, Cong Liu, and Ji-Rong Wen. 2022. Ji-	2170
2122	Tianlu Wang, and Luke Zettlemoyer. 2022b. OPT:	uzhang: A chinese pre-trained language model for	2171
2123	open pre-trained transformer language models . <i>ArXiv</i>	mathematical problem understanding . In <i>KDD ’22:</i>	2172
2124	<i>preprint</i> , abs/2205.01068 .	<i>The 28th ACM SIGKDD Conference on Knowledge</i>	2173
2125	Tianhua Zhang, Jiaxin Ge, Hongyin Luo, Yung-Sung	<i>Discovery and Data Mining, Washington, DC, USA,</i>	2174
2126	Chuang, Mingye Gao, Yuan Gong, Xixin Wu, Yoon	<i>August 14 - 18, 2022</i> , pages 4571–4581. ACM.	2175
2127	Kim, Helen Meng, and James Glass. 2023d. Natural	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	2176
2128	language embedded programs for hybrid language	Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen	2177
2129	symbolic reasoning . <i>ArXiv preprint</i> , abs/2309.10814 .	Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen	2178
		Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang,	2179
		Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu,	2180
		Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023b.	2181
		A survey of large language models . <i>ArXiv preprint</i> ,	2182
		abs/2303.18223 .	2183
		Xufeng Zhao, Mengdi Li, Wenhao Lu, Cornelius Weber,	2184
		Jae Hee Lee, Kun Chu, and Stefan Wermter. 2023c.	2185

2186	Enhancing zero-shot chain-of-thought reasoning in	Yuchen Zhuang, Xiang Chen, Tong Yu, Saayan Mitra,	2243
2187	large language models through logic. <i>ArXiv preprint,</i>	Victor Bursztyn, Ryan A. Rossi, Somdeb Sarkhel,	2244
2188	abs/2309.13339.	and Chao Zhang. 2023. Toolchain* : Efficient action	2245
2189	Chuanyang Zheng, Zhengying Liu, Enze Xie, Zhenguo	space navigation in large language models with a*	2246
2190	Li, and Yu Li. 2023a. Progressive-hint prompting	search.	2247
2191	improves reasoning in large language models. <i>ArXiv</i>	Anni Zou, Zhuosheng Zhang, Hai Zhao, and Xian-	2248
2192	<i>preprint</i> , abs/2304.09797.	gru Tang. 2023. Meta-cot: Generalizable chain-of-	2249
2193	Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and	thought prompting in mixed-task scenarios with large	2250
2194	Sibei Yang. 2023b. DDCot: Duty-distinct chain-of-	language models. <i>ArXiv preprint</i> , abs/2310.06692.	2251
2195	thought prompting for multimodal reasoning in lan-		
2196	guage models. In <i>Thirty-seventh Conference on Neu-</i>		
2197	<i>ral Information Processing Systems, NeurIPS 2023.</i>		
2198	Huaxiu Steven Zheng, Swaroop Mishra, Xinyun Chen,		
2199	Heng-Tze Cheng, Ed H Chi, Quoc V Le, and Denny		
2200	Zhou. 2023c. Take a step back: Evoking reason-		
2201	ing via abstraction in large language models. <i>ArXiv</i>		
2202	<i>preprint</i> , abs/2310.06117.		
2203	Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman,		
2204	Haohan Wang, and Yu-Xiong Wang. 2023a. Lan-		
2205	guage agent tree search unifies reasoning acting and		
2206	planning in language models.		
2207	Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth.		
2208	2019. “going on a vacation” takes longer than “go-		
2209	ing for a walk” : A study of temporal commonsense		
2210	understanding. In <i>Proceedings of the 2019 Confer-</i>		
2211	<i>ence on Empirical Methods in Natural Language Pro-</i>		
2212	<i>cessing and the 9th International Joint Conference</i>		
2213	<i>on Natural Language Processing (EMNLP-IJCNLP),</i>		
2214	pages 3363–3369, Hong Kong, China. Association		
2215	for Computational Linguistics.		
2216	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei,		
2217	Nathan Scales, Xuezhi Wang, Dale Schuurmans,		
2218	Claire Cui, Olivier Bousquet, Quoc V. Le, and Ed H.		
2219	Chi. 2023b. Least-to-most prompting enables com-		
2220	plex reasoning in large language models. In <i>The</i>		
2221	<i>Eleventh International Conference on Learning Rep-</i>		
2222	<i>resentations, ICLR 2023, Kigali, Rwanda, May 1-5,</i>		
2223	2023. OpenReview.net.		
2224	Yuxiang Zhou, Jiazheng Li, Yanzheng Xiang, Hanqi		
2225	Yan, Lin Gui, and Yulan He. 2023c. The mystery		
2226	and fascination of llms: A comprehensive survey on		
2227	the interpretation and analysis of emergent abilities.		
2228	<i>ArXiv preprint</i> , abs/2311.00237.		
2229	Zhehua Zhou, Jiayang Song, Kunpeng Yao, Zhan Shu,		
2230	and Lei Ma. 2023d. Isr-llm: Iterative self-refined		
2231	large language model for long-horizon sequential		
2232	task planning.		
2233	Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao		
2234	Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-		
2235	Seng Chua. 2021. TAT-QA: A question answering		
2236	benchmark on a hybrid of tabular and textual con-		
2237	tent in finance. In <i>Proceedings of the 59th Annual</i>		
2238	<i>Meeting of the Association for Computational Lin-</i>		
2239	<i>guistics and the 11th International Joint Conference</i>		
2240	<i>on Natural Language Processing (Volume 1: Long</i>		
2241	<i>Papers)</i> , pages 3277–3287, Online. Association for		
2242	Computational Linguistics.		

A Appendix

A.1 Related Survey

Zhao et al. (2023b) primarily focuses on the development of contemporary LLMs, while Qiu et al. (2020) surveys about early PLMs. Some works discuss reasoning in specific domains, such as mathematical reasoning (Lu et al., 2023c), common-sense reasoning (Talmor et al., 2019), and logical reasoning (Yang et al., 2023c). Huang et al. (2023c); Zhang et al. (2023f) conducts an investigation into potential hallucination phenomena in LLM’s reasoning. Dong et al. (2023) discusses in-context learning techniques in the era of LLMs, and Yu et al. (2023a) conducts a macroscopic investigation into natural language reasoning. Liu et al. (2023d) mainly discusses prompt tuning, while Qiao et al. (2023); Yu et al. (2023d) are more concentrated on prompt engineering and strategies.

Distinct from the above-mentioned surveys, this paper focuses on generalized chain-of-thought reasoning in the era of LLMs. This is the first systematic investigation into XoT reasoning, and we hope our work can serve as an overview to facilitate future research.

A.2 Further Discussion

Open Question: Does CoT ability originate from code data pre-training? This is a pending question, initially summarized by Fu and Khot (2022) and widely circulated in the research community. In the early stages, LLMs like GPT3 (Brown et al., 2020) (davinci) and OPT (Zhang et al., 2022b) usually do not possess CoT capabilities and they do not use or only incorporate a small amount of code data (not specialized) during pre-training. Recent models often incorporate specialized code data during pre-training, such as GPT-3.5, LLaMA2 (Touvron et al., 2023b) (with approximately 8% of code data during pre-training) and they all possess strong CoT capabilities. Additionally, Gao et al. (2023); Chen et al. (2022a) have found that the use of programming language form rationales can significantly enhance the model’s performance on complex reasoning tasks. Various indications point towards the source of CoT abilities lying in code data during pre-training. However, as of now, there is no work that has reached a definite conclusion on this opinion, which necessitates further in-depth exploration in future research.

Open Question: How to provide precise feedback on model’s reasoning or decisions? When dealing with multi-step reasoning or decision-making tasks, errors often occur in intermediate steps, and if these errors are not corrected promptly, they may lead to cascading errors. Currently, the primary methods for obtaining feedback include feedback from model itself (Madaan et al., 2023; Shinn et al., 2023), feedback from other models (Paul et al., 2023), feedback from the external environment (Nathani et al., 2023; Gou et al., 2023a), and feedback based on reinforcement learning (Uesato et al., 2022; Lightman et al., 2023; Ma et al., 2023). However, these methods have inherent issues. (1) How dependable is the feedback generated by the model itself? (2) Is there a fundamental distinction between feedback from other models and self-feedback? (3) Does the feedback quality still remain constrained by the model’s capability boundaries? (4) How is external feedback for various scenarios pre-defined, and how can this be expanded to different scenarios? (5) How to obtain an effective reward model?

In summary, there is currently no fully satisfying feedback approach and more research attention is needed on how to accurately obtain feedback signals from the model’s intermediate processes.

Discussion: Towards (early) AGI AGI has been the long-standing ultimate aspiration in the realm of artificial intelligence.

Integration of reasoning and world interaction. With robust language comprehension capabilities, LLMs can engage with the external world through text-based interactions using plugins (tools, API, etc) (Schick et al., 2023; Shen et al., 2023a; Qin et al., 2023b). Combining powerful reasoning capabilities, LLMs have made significant strides in various planning and decision-making tasks (Shinn et al., 2023; Yao et al., 2023b; Zhuang et al., 2023), catalyzing research on LLM-based autonomous agents (Wang et al., 2023h; Xi et al., 2023).

LLM acts as the Brain (Controller). In contrast to traditional AI, which concentrates on specific tasks, AGI seeks the ability to understand general tasks (Devlin et al., 2019; Dosovitskiy et al., 2021), covering a widespread spectrum. Within LLM-based AI, the LLM typically serves as the brain (or central controller), handling reasoning, planning and decision-making, while delegating specific execution to dedicated modules (tools, weak AI, etc.) (Shen et al., 2023a; Yang et al.,

2023a). LLM-based AI has already diverged significantly from weak AI and is progressing towards human cognition and thinking.

While some studies suggest that LLMs represent an early manifestation of AGI (Bubeck et al., 2023; Jack, 2023), there are also scholars who contend that LLMs may not progress into AGI due to factors such as auto-regressive modeling and limited memory. As of now, there is still intense debate on whether LLMs can evolve into AGI. But regardless, LLM-based AI has embarked on a distinctly different path from traditional AI, evolving towards a more generalized direction.

A.3 Early Attempts and Efforts in Specific Domains

In this section, we list the early attempts of XoT reasoning and efforts focused on specific domains.

Before the concept of CoT was introduced (Wei et al., 2022b), some efforts were made to enhance reasoning performance through the use of rationales (Marasovic et al., 2022; Rajani et al., 2019a,b; Dua et al., 2020). After that, certain work has empirically demonstrated the effectiveness of chain-of-thought prompting (Lampinen et al., 2022; Ye and Durrett, 2022; Arora et al., 2023) and Shi et al. (2023) explores multi-lingual CoT reasoning. Other work focuses on specific domains, such as machine translation (He et al., 2023b), sentiment analysis (Fei et al., 2023), sentence embeddings (Zhang et al., 2023a), summarization (Wang et al., 2023n), arithmetic (Lee and Kim, 2023), and tabular reasoning (Chen, 2023; Jin and Lu, 2023), etc. Katz et al. (2022); Zhang et al. (2022a) provide benchmarks and resources. Besides, some research utilizes specific pre-training to enhance certain capabilities, such as mathematical reasoning (Lewkowycz et al., 2022; Zhao et al., 2022).

A.4 Empirical Results

We statistic the performance of various XoT methods in mathematics, common sense, and symbolic reasoning, as shown in Table 2. We primarily focus on the performance of GPT series models and the results are mainly from corresponding papers (some results are used as baselines in other papers). It is worth noting that due to variations in model versions and experimental setups, even the methods with the same backbone may not be fairly comparable on the same dataset. Therefore, this table only provides trends and empirical insights.

B Details of Benchmarks

B.1 Mathematical Reasoning

Mathematical reasoning is often used to measure the reasoning power of a model. Early benchmarks contain simple arithmetic operations (Hosseini et al., 2014; Koncel-Kedziorski et al., 2015; Roy and Roth, 2015; Koncel-Kedziorski et al., 2016). Ling et al. (2017) labels the reasoning process in natural language form, and Amini et al. (2019) builds on AQUA by labeling the reasoning process in program form. Later benchmarks (Miao et al., 2020; Patel et al., 2021; Cobbe et al., 2021; Gao et al., 2023) contain more complex and diverse questions. (Zhu et al., 2021; Chen et al., 2021, 2022b) require reasoning based on the table content. There are also general benchmarks (Hendrycks et al., 2021b; Mishra et al., 2022a,b) and reading comprehension form benchmarks (Dua et al., 2019; Chen et al., 2023a).

B.2 Commonsense Reasoning

Commonsense reasoning is the process of making inferences, judgments, and understandings based on knowledge that is generally known and commonly perceived in the everyday world. How to acquire and understand commonsense knowledge is a major impediment to models facing commonsense reasoning. Many benchmarks and tasks are proposed focusing on commonsense understanding (Talmor et al., 2019, 2021; Bhakthavatsalam et al., 2021; Mihaylov et al., 2018; Geva et al., 2021; Huang et al., 2019; Bisk et al., 2020), event temporal commonsense reasoning (Rashkin et al., 2018; Zhou et al., 2019), and commonsense verification (Wang et al., 2019).

B.3 Symbolic Reasoning

Symbolic reasoning here refers specifically to the simulation of some simple operations, which are simple for humans yet challenging for LLMs. Last letter concatenation, coin flip, and reverse list (Wei et al., 2022b) are the most commonly used symbolic reasoning tasks. In addition, the collaborative benchmark BigBench (Srivastava et al., 2022) and BigBench-Hard (Suzgun et al., 2023) also contain several symbolic reasoning datasets, such as state tracking and object counting.

B.4 Logical Reasoning

Logical reasoning is divided into deductive reasoning, inductive reasoning, and abductive reason-

Task	Dataset	Size	Input	Output	Rationale	Description
Mathematical Reasoning	AddSub (Hosseini et al., 2014)	395	Question	Number	Equation	Simple arithmetic
	SingleEq (Koncel-Kedziorski et al., 2015)	508	Question	Number	Equation	Simple arithmetic
	MultiArith (Roy and Roth, 2015)	600	Question	Number	Equation	Simple arithmetic
	MAWPS (Koncel-Kedziorski et al., 2016)	3320	Question	Number	Equation	Simple arithmetic
	AQUA-RAT (Ling et al., 2017)	100,000	Question	Option	Natural Language	Math reasoning with NL rationale
	ASDiv (Miao et al., 2020)	2305	Question	Number	Equation	Multi-step math reasoning
	SVAMP (Patel et al., 2021)	1,000	Question	Number	Equation	Multi-step math reasoning
	GSM8K (Cobbe et al., 2021)	8,792	Question	Number	Natural Language	Natural Language
	GSM-Hard (Gao et al., 2023)	936	Question	Number	Natural Language	GSM8K with larger number
	MathQA (Amini et al., 2019)	37,297	Question	Number	Operation	Annotated based on AQUA
	DROP (Dua et al., 2019)	96,567	Question+Passage	Number+Span	Equation	Reading comprehension form
	TheoremQA (Chen et al., 2023a)	800	Question+Theorem	Number	✗	Answer based on theorems
	TAT-QA (Zhu et al., 2021)	16,552	Question+Table+Text	Number+Span	Operation	Answer based on tables
	FinQA (Chen et al., 2021)	8,281	Question+Table+Text	Number	Operation	Answer based on tables
	ConvFinQA (Chen et al., 2022b)	3892	Question+Table+Dialog	Number	Operation	Multi-turn dialogs
	MATH (Hendrycks et al., 2021b)	12500	Question	Number	Natural Language	Challenging competition math problems
Commonsense Reasoning	NumGLUE (Mishra et al., 2022b)	101,835	Question+Text	Number+Span	✗	Multi-task benchmark
	LILA (Mishra et al., 2022a)	133,815	Question+Text	Free-form	Program	Multi-task benchmark
	ARC (Bhaskaravatsalam et al., 2021)	7787	Question	Option	✗	From science exam
	OpenBookQA (Mihaylov et al., 2018)	5,957	Question+Context	Option	✗	Open-book knowledge
	PIQA (Bisk et al., 2020)	21000	Goal+Solution	Option	✗	Physical commonsense knowledge
	CommonsenseQA (Talmor et al., 2019)	12247	Question	Option	✗	Derived from ConceptNet
	CommonsenseQA 2.0 (Talmor et al., 2021)	14343	Question	Yes/No	✗	Gaming annotation with high quality
	Event2Mind (Rashkin et al., 2018)	25000	Event	Intent+Reaction	✗	Intension commonsense reasoning
	McTaco (Zhou et al., 2019)	13225	Question	Option	✗	Event temporal commonsense reasoning
	CosmosQA (Huang et al., 2019)	35588	Question+Paragraph	Option	✗	Narrative commonsense reasoning
Symbolic Reasoning	ComValidation (Wang et al., 2019)	11997	Statement	Option	✗	Commonsense verification
	ComExplanation (Wang et al., 2019)	11997	Statement	Option/Free-form	✗	Commonsense explanation
	StrategyQA (Geva et al., 2021)	2,780	Question	Yes/No	✗	Multi-hop commonsense reasoning
	Last Letter Concat. (Wei et al., 2022b)	-	Words	Letters	✗	Rule-based
	Coin Flip (Wei et al., 2022b)	-	Statement	Yes/No	✗	Rule-based
Logical Reasoning	Reverse List (Wei et al., 2022b)	-	List	Reversed List	✗	Rule-based
	BigBench (Srivastava et al., 2022)	-	-	-	✗	Contains multiple symbolic reasoning datasets
	BigBench-Hard (Suzgun et al., 2023)	-	-	-	✗	Contains multiple symbolic reasoning datasets
	ReClor (Yu et al., 2020)	6,138	Question+Context	Option	✗	Questions from GMAT and LSAT
	LogiQA (Liu et al., 2020)	8,678	Question+Paragraph	Option	✗	Questions from China Civil Service Exam
	ProofWriter (Tafjord et al., 2021)	20192	Question+Rule	Answer+Proof	Entailment Tree	Reasoning process generation
	FOLIO (Han et al., 2022)	1435	Conclusion+Premise	Yes/No	✗	First-order logic
Multimodal Reasoning	DEER (Yang et al., 2022)	1,200	Fact	Rule	✗	Inductive reasoning
	PrOntoQA (Saparov and He, 2023)	-	Question+Context	Yes/No+Process	First-Order Logic	Deductive reasoning
	VCR (Zellers et al., 2019)	264,720	Question+Image	Option	Natural Language	Visual commonsense reasoning
	VisualCOMET (Park et al., 2020)	1,465,704	Image+Event	Action+Intent	✗	Visual commonsense reasoning
	PMR (Dong et al., 2022)	15,360	Image+Background	Option	✗	Premise-based multi-modal reasoning
	ScienceQA (Lu et al., 2022)	21,208	Q+Image+Context	Option	Natural Language	Multi-modal reasoning with NL rationales
	VLEP (Lei et al., 2020)	28,726	Premise+Video	Option	✗	Video event prediction
	CLEVRER (Yi et al., 2020)	305,280	Question+Video	Option/Free-form	Program	Video temporal and causal reasoning
	STAR (Wu et al., 2021)	600,000	Question+Video	Option	✗	Video situated reasoning
	NEXT-QA (Xiao et al., 2021)	47,692	Question+Video	Option	✗	Video temporal, causal, commonsense reasoning
	Causal-VidQA (Li et al., 2022a)	107,600	Question+Video	Free-form	Natural Language	Video causal and commonsense reasoning
	News-KVQA (Gupta and Gupta, 2022)	1,041,352	Q+V+KG	Option	✗	Video reasoning with external knowledge

Table 1: An overview of benchmarks and tasks on reasoning.

ing (Yu et al., 2023a). Deductive reasoning derives conclusions from general premises (Liu et al., 2020; Yu et al., 2020; Tafjord et al., 2021; Han et al., 2022). Inductive reasoning derives general conclusions from special cases (Yang et al., 2022). Abductive reasoning gives rational explanations for observed phenomena (Saparov and He, 2023).

B.5 Multi-modal Reasoning

In the real world, reasoning also involves information in modalities other than text, with visual modalities being the most prevalent. To this end, many benchmarks for visual multi-modal reasoning are proposed (Zellers et al., 2019; Park et al., 2020; Dong et al., 2022; Lu et al., 2022), and among them, ScienceQA (Lu et al., 2022) annotates reasoning process and is the most commonly used visual multi-modal reasoning benchmark. Video multi-modal reasoning (Lei et al., 2020; Yi et al., 2020; Wu et al., 2021; Xiao et al., 2021; Li et al., 2022a; Gupta and Gupta, 2022) is more challenging as it introduces additional temporal information compared to visual multi-modal reasoning.

B.6 Comprehensive Benchmarks

Apart from the aforementioned individual datasets, there are also some comprehensive evaluation benchmarks. Some works aim to provide a holistic evaluation of the general reasoning capabilities (Srivastava et al., 2022; Suzgun et al., 2023; Hendrycks et al., 2021a; Huang et al., 2023e; Liang et al., 2022). In addition, there are also some multi-task benchmarks that focus on specific reasoning abilities, such as logical reasoning (Luo et al., 2023; Liu et al., 2023b) and temporal reasoning (Chu et al., 2023; Wang and Zhao, 2023).

B.7 Evaluation Metrics

Accuracy Accuracy is used to assess a model’s ability on classification tasks and is commonly used for multi-choice (Ling et al., 2017; Mihaylov et al., 2018; Liu et al., 2020; Lu et al., 2022) and yes/no (Talmor et al., 2021; Geva et al., 2021; Han et al., 2022) tasks.

$$\text{Accuracy} = \frac{N_{\text{correct}}}{N_{\text{total}}} \quad (7)$$

Method	Setting	Backbone	Mathematical				Commonsense		Symbolic	
			GSM8K	SVAMP	Asdiv	AQuA	CSQA	StrategyQA	LastLetterConcat	CoinFlip
I-O Prompting (Brown et al., 2020)	fewshot	text-davinci-002	19.7	69.9	74	29.5	79.5	65.9	5.8	49.0
Fewshot CoT (Wei et al., 2022b)	fewshot	text-davinci-002	63.1	76.4	80.4	45.3	73.5	65.4	77.5	99.6
PoT (Chen et al., 2022a)	fewshot	text-davinci-002	80	89.1	-	58.6	-	-	-	-
Complex CoT (Fu et al., 2023a)	fewshot	text-davinci-002	72.6	-	-	-	-	77	-	-
Automate CoT (Shum et al., 2023)	fewshot	text-davinci-002	49.7	73.3	74.2	37.9	76.1	67.9	58.9	-
Fewshot CoT (Wei et al., 2022b)	fewshot	text-davinci-003	16.83	69.06	-	29.13	-	-	-	-
PHP (Zheng et al., 2023a)	fewshot	text-davinci-003	79	84.7	-	58.6	-	-	-	-
Self-consistency (Wang et al., 2023m)	fewshot	text-davinci-003	67.93	83.11	-	55.12	-	-	-	-
Active Prompt (Diao et al., 2023)	fewshot	text-davinci-003	65.6	80.5	79.8	48	78.9	74.2	71.2	-
Synthetic Prompt (Shao et al., 2023b)	fewshot	text-davinci-003	73.9	81.8	80.7	-	-	-	-	-
FOBAR (Jiang et al., 2023b)	fewshot	text-davinci-003	79.5	86	-	58.66	-	-	-	-
Boosted Prompting (Pitis et al., 2023)	fewshot	text-davinci-003	71.6	-	-	55.1	-	-	-	-
Fewshot CoT (Wei et al., 2022b)	fewshot	code-davinci-002	60.1	75.8	80.1	39.8	79	73.4	70.4	99
Self-Consistency (Wang et al., 2023m)	fewshot	code-davinci-002	78	86.8	87.8	52	81.5	79.8	73.4	99.5
PAL (Gao et al., 2023)	fewshot	code-davinci-002	72	79.4	79.6	-	-	-	-	-
Resprompt (Jiang et al., 2023a)	fewshot	code-davinci-002	66.6	-	-	45.3	-	-	-	-
DIVERSE (Li et al., 2022c)	fewshot	code-davinci-002	82.3	87	88.7	-	79.9	78.6	-	-
Least-to-Most (Zhou et al., 2023b)	fewshot	code-davinci-002	68.01	-	-	-	-	-	94	-
Boosted Prompting (Pitis et al., 2023)	fewshot	code-davinci-002	83.3	88.6	-	61.7	-	-	-	-
Fewshot CoT (Wei et al., 2022b)	fewshot	gpt-3.5-turbo	76.5	81.9	-	54.3	78	63.7	73.2	99
Self-consistency (Wang et al., 2023m)	fewshot	gpt-3.5-turbo	81.9	86.4	-	62.6	-	-	-	-
MetaCoT (Zou et al., 2023)	fewshot	gpt-3.5-turbo	75.1	88.6	-	54.7	72.4	64.5	77.2	100
Verify CoT (Ling et al., 2023)	fewshot	gpt-3.5-turbo	86	-	-	69.5	-	-	92.6	-
Active Prompting (Diao et al., 2023)	fewshot	gpt-3.5-turbo	81.8	82.5	87.9	55.3	-	-	-	-
RCoT (Xue et al., 2023)	fewshot	gpt-3.5-turbo	84.6	84.9	89.3	57.1	-	-	-	-
FOBAR (Jiang et al., 2023b)	fewshot	gpt-3.5-turbo	87.4	87.4	-	57.5	-	-	-	-
Memory-of-Thought (Li and Qiu, 2023)	fewshot	gpt-3.5-turbo	-	-	-	54.1	-	-	-	-
Adaptive-consistency (Aggarwal et al., 2023)	fewshot	gpt-3.5-turbo	82.7	85	83	-	-	67.9	-	-
Boosted Prompting (Pitis et al., 2023)	fewshot	gpt-3.5-turbo	87.1	-	-	72.8	-	-	-	-
Zeroshot CoT (Kojima et al., 2022)	zeroshot	text-davinci-002	40.5	63.7	-	31.9	64	52.3	57.6	87.8
PoT (Chen et al., 2022a)	zeroshot	text-davinci-002	57	70.8	-	43.9	-	-	-	-
AutoCoT (Zhang et al., 2023h)	zeroshot	text-davinci-002	47.9	69.5	-	36.5	74.4	65.4	59.7	99.9
COSP (Aggarwal et al., 2023)	zeroshot	code-davinci-001	8.7	-	-	-	55.4	52.8	-	-
Plan-and-Solve (Wang et al., 2023i)	zeroshot	text-davinci-003	58.2	72	-	42.5	65.2	63.8	64.8	96.8
Agent-Instruct (Crispino et al., 2023)	zeroshot	gpt-3.5-turbo	73.4	80.8	-	57.9	74.1	69	99.8	95.2
Self-Refine (Madaan et al., 2023)	zeroshot	gpt-3.5-turbo	64.1	-	-	-	-	-	-	-
RCoT (Xue et al., 2023)	zeroshot	gpt-3.5-turbo	82	79.6	86	55.5	-	-	-	-

Table 2: The performance of various XoT methods in commonly used mathematical, commonsense and symbolic reasoning benchmarks. It is worth noting that, due to variations in the experimental setups of different methods, their performances are not directly comparable. The table is used to provide an overall empirical insight.

EM and F1 EM and F1 are metrics used to evaluate free form (Mishra et al., 2022a; Wang et al., 2019; Yi et al., 2020) and span extraction (Dua et al., 2019; Zhu et al., 2021; Mishra et al., 2022b) tasks. Both are calculated at the token level.

$$F1 = \frac{2 \cdot P \cdot R}{P + R} \quad (8)$$

$$EM = \frac{\sum \mathbb{I}[A = A']}{N_{\text{total}}} \quad (9)$$

where P and R stand for precision and recall, and EM calculates the proportion of predictions and answers that are exactly the same.

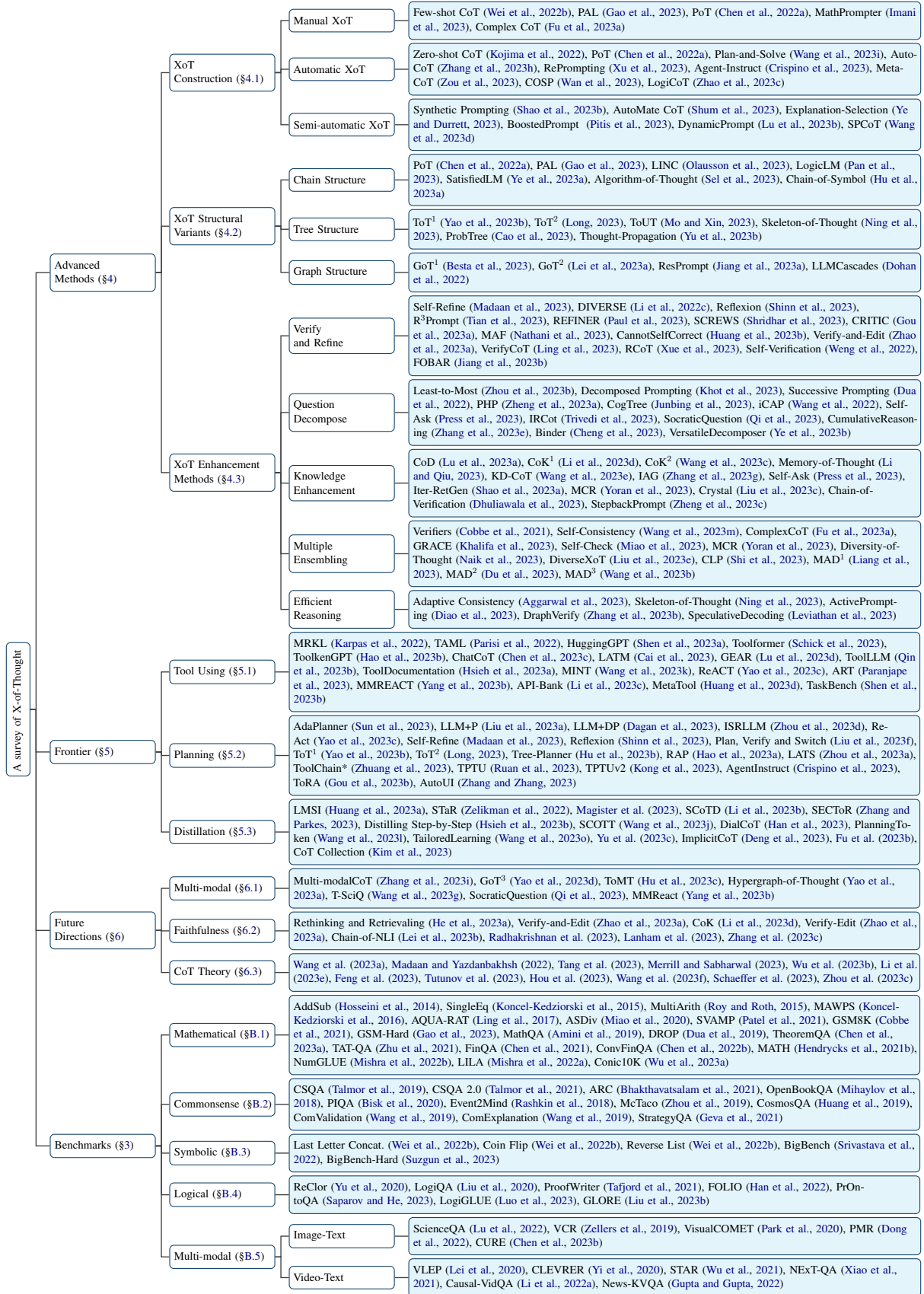


Figure 8: Taxonomy of Advanced Methods, Frontiers, Future Directions, and Benchmarks (Full Edition).