


From System 1 to System 2: A Survey of Reasoning Large Language Models

Zhong-Zhi Li^{*}, Duzhen Zhang^{*}, Ming-Liang Zhang[§], Jiaxin Zhang[§], Zengyan Liu[§], Yuxuan Yao[§],
Haotian Xu, Junhao Zheng, Pei-Jie Wang, Xiuyi Chen, Yingying Zhang, Fei Yin, Jiahua Dong,
Zhijiang Guo[†], Le Song[†], Cheng-Lin Liu[†] , *Fellow, IEEE*

Abstract—Achieving human-level intelligence requires refining the transition from the fast, intuitive System 1 to the slower, more deliberate System 2 reasoning. While System 1 excels in quick, heuristic decisions, System 2 relies on logical reasoning for more accurate judgments and reduced biases. Foundational Large Language Models (LLMs) excel at fast decision-making but lack the depth for complex reasoning, as they have not yet fully embraced the step-by-step analysis characteristic of true System 2 thinking. Recently, reasoning LLMs like OpenAI’s o1/o3 and DeepSeek’s R1 have demonstrated expert-level performance in fields such as mathematics and coding, closely mimicking the deliberate reasoning of System 2 and showcasing human-like cognitive abilities. This survey begins with a brief overview of the progress in foundational LLMs and the early development of System 2 technologies, exploring how their combination has paved the way for reasoning LLMs. Next, we discuss how to construct reasoning LLMs, analyzing their features, the core methods enabling advanced reasoning, and the evolution of various reasoning LLMs. Additionally, we provide an overview of reasoning benchmarks, offering an in-depth comparison of the performance of representative reasoning LLMs. Finally, we explore promising directions for advancing reasoning LLMs and maintain a real-time GitHub Repository to track the latest developments. We hope this survey will serve as a valuable resource to inspire innovation and drive progress in this rapidly evolving field.

Index Terms—Slow-thinking, Large Language Models, Human-like Reasoning, Decision Making in AI, AGI



1 INTRODUCTION

“Don’t teach. Incentivize.”

—Hyung Won Chung, OpenAI

ACHIEVING human-level intelligence requires refining the transition from *System 1* to *System 2* reasoning [1]–[5]. Dual-system theory suggests that human cognition operates through two modes: *System 1*, which is fast, automatic, and intuitive, enabling quick decisions with minimal effort, and *System 2*, which is slower, more analytical, and deliberate [6], [7]. While *System 1* is efficient for routine

tasks, it is prone to cognitive biases, especially in complex or uncertain situations, leading to judgment errors. In contrast, *System 2* relies on logical reasoning and systematic thinking, resulting in more accurate and rational decisions [8]–[11]. By mitigating the biases of *System 1*, *System 2* provides a more refined approach to problem-solving [12]–[15].

The development of foundational Large Language Models (LLMs)¹ has marked a major milestone in Artificial Intelligence (AI). Models such as GPT-4o [16] and DeepSeek-v3 [17] have demonstrated impressive capabilities in text generation, language translation, and a variety of perception tasks [18]–[28]. These models, trained on extensive datasets and utilizing advanced algorithms, excel in understanding and generating human-like responses. However, despite their impressive achievements, foundational LLMs operate in a manner similar to *System 1* reasoning, relying on fast, heuristic-driven decision-making. While they perform exceptionally well in providing rapid responses, they often fall short in scenarios requiring deep, logical analysis and precision in complex reasoning tasks. This limitation becomes especially clear in situations involving intricate problem-solving, logical analysis, or nuanced understanding, where these models do not yet match human cognitive abilities.

In contrast, reasoning LLMs represent a significant advancement in the evolution of language models. Models

Version: v1 (major update on February 23, 2025)

^{*}Core contribution. [§]Significant contribution. [†]Corresponding author.

Duzhen Zhang, Jiahua Dong, and Le Song are with the Mohamed bin Zayed University of Artificial Intelligence, Abu Dhabi, UAE (E-mail: bladedancer957@gmail.com; dongjiahua1995@gmail.com; le.song@mbzuai.ac.ae).

Zhong-Zhi Li, Pei-Jie Wang, Xiuyi Chen, Fei Yin, and Cheng-Lin Liu are with the Institute of Automation, Chinese Academy of Sciences, Beijing, China (E-mail: lizhongzhi2022@ia.ac.cn; wangpeijie2023@ia.ac.cn; hugheren.chan@gmail.com; fyin@nlpr.ia.ac.cn; liucl@nlpr.ia.ac.cn).

Ming-Liang Zhang is with the AiShiWeiLai AI Research, Beijing, China (E-mail: zhangmingliang@yuaiweiwei.com).

Zengyan Liu, Yuxuan Yao, and Zhijiang Guo is with City University of Hong Kong and the Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China (E-mail: zengyalu2-c@my.cityu.edu.hk; yuxuan Yao3-c@my.cityu.edu.hk; zhijiangguo@hkust-gz.edu.cn).

Jiaxin Zhang is with the University of Strathclyde, Glasgow, UK (E-mail: jiaxin.zhang@strath.ac.uk).

Haotian Xu is with the Xiaohongshu Inc, Beijing, China (E-mail: xuhao-tian@xiaohongshu.com).

Yingying Zhang is with the East China Normal University, Shanghai, China (E-mail: yyzhang@fem.ecnu.edu.cn).

Junhao Zheng is with the South China University of Technology, Guangzhou, China (E-mail: junhaozheng47@outlook.com).

1. In this paper, “reasoning” refers to answering questions involving complex, multi-step processes with intermediate steps. **Foundational LLMs:** LLMs with basic reasoning abilities, handling simple or single-step tasks. **Reasoning LLMs:** LLMs that excel in complex tasks like coding and mathematical proofs, incorporating a “thinking” process—tasks that foundational LLMs struggle with.

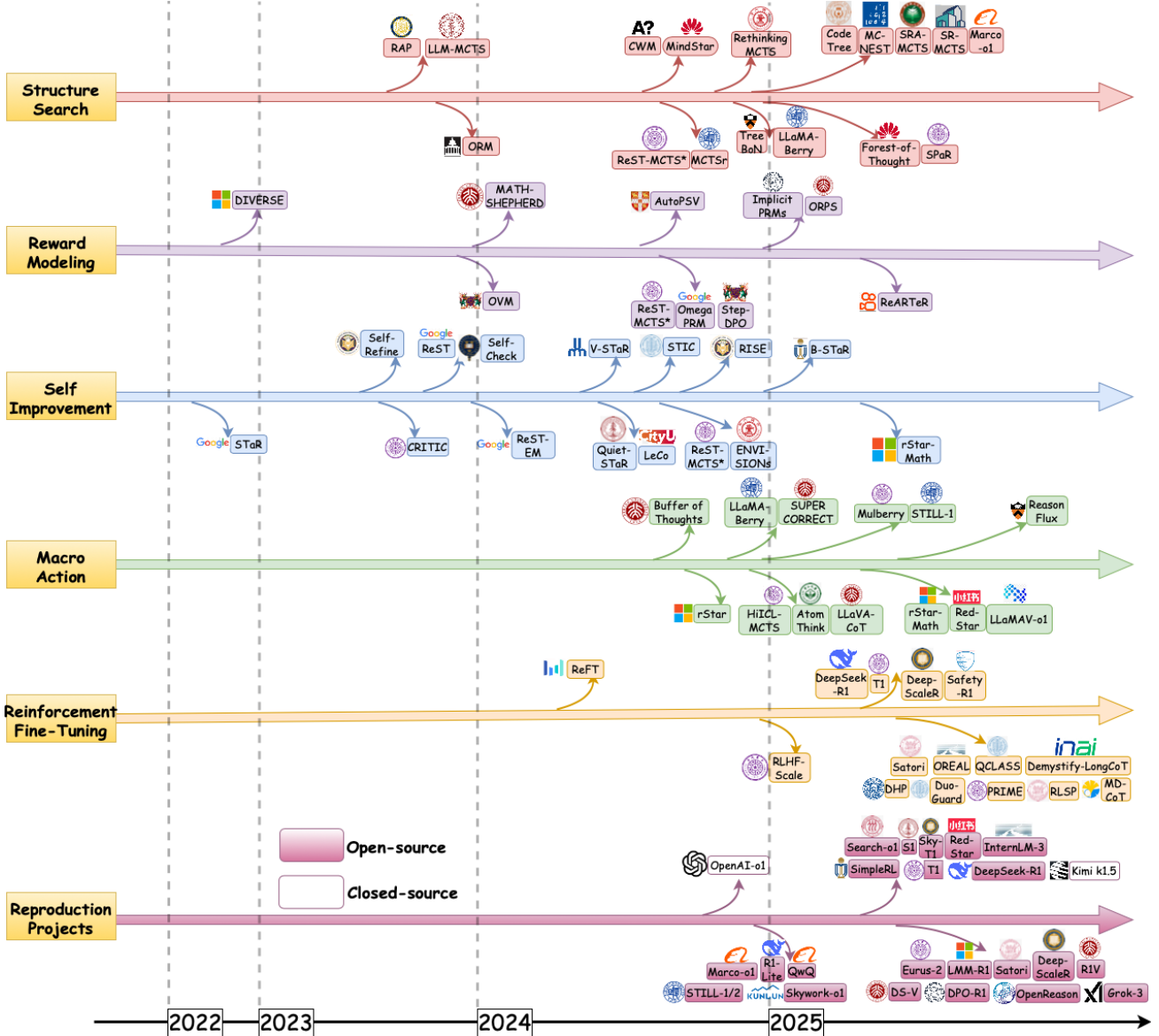


Fig. 1. The recent timeline of reasoning LLMs, covering core methods and the release of open-source and closed-source reproduction projects.

like OpenAI’s o1/o3 [29], [30] and DeepSeek’s R1 [31] are designed to emulate the slower, more deliberate reasoning associated with *System 2* thinking. Unlike foundational LLMs, reasoning LLMs are equipped with mechanisms for processing information step-by-step, allowing them to make more accurate and rational decisions. This shift from fast-thinking, intuitive processes to more methodical, reasoning-driven models enables reasoning LLMs to tackle complex tasks, such as advanced mathematics [32]–[37], logical reasoning [38]–[44], and multimodal reasoning [45]–[47], with expert-level performance, exhibiting human-like cognitive abilities. As a result, reasoning LLMs are increasingly seen as capable of achieving the deep, logical thinking needed for tasks that were once considered beyond AI’s reach. The recent timeline of reasoning LLMs is presented in Figure 1.

1.1 Structure of the Survey

This survey offers a comprehensive overview of the key concepts, methods, and challenges involved in the development

of reasoning LLMs. As illustrated in Figure 2, this survey is organized as follows:

- 1) Section 2 offers a concise overview of the progress in foundational LLMs (Section 2.1) and the early development of key *System 2* technologies, including symbolic logic systems (Section 2.2), Monte Carlo Tree Search (MCTS) (Section 2.3), and Reinforcement Learning (RL) (Section 2.4), highlighting how their combination has paved the way for reasoning LLMs.
- 2) Section 3 introduces reasoning LLMs and outlines their construction process. Specifically, Section 3.1 presents the characteristics of reasoning LLMs from two perspectives: output behavior (Section 3.1.1) and training dynamics (Section 3.1.2), emphasizing their differences from foundational LLMs. Section 3.2 identifies the core methods necessary for achieving advanced reasoning capabilities, focusing on five aspects: Structure Search (Section 3.2.1), Reward Modeling (Section 3.2.2), Self Improvement (Section 3.2.3), Macro Action (Section

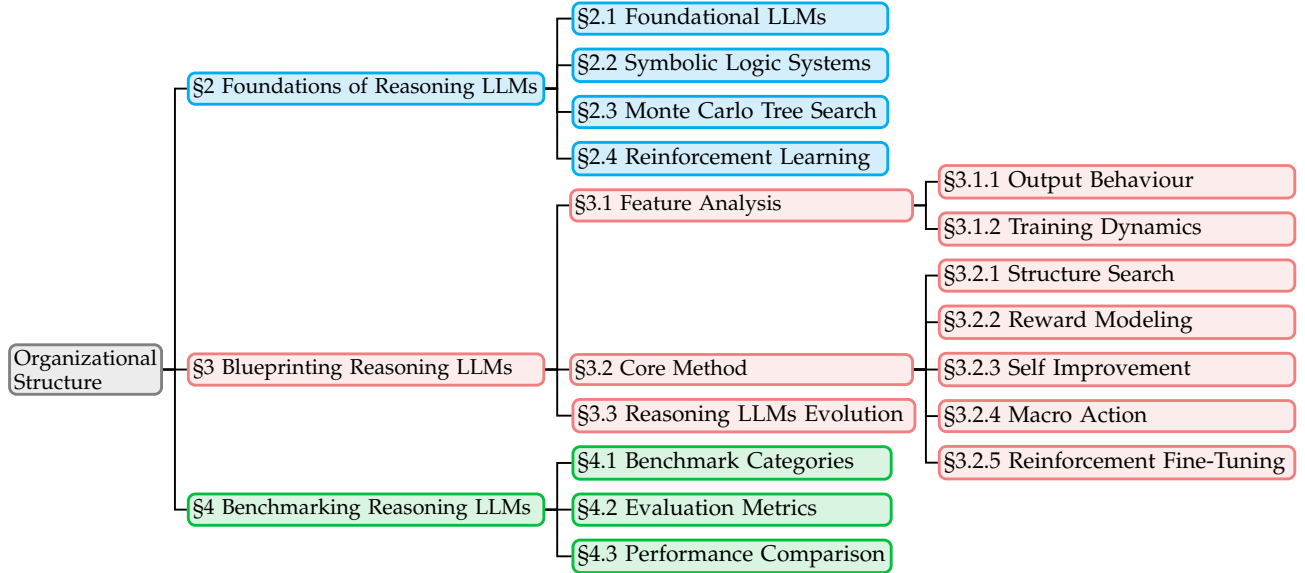


Fig. 2. The primary organizational structure of the survey.

3.2.4), and Reinforcement Fine-Tuning (Section 3.2.5). Each section delves into the specific characteristics of these methods and introduces representative reasoning LLMs for each approach. Section 3.3 traces the evolutionary stages of reasoning LLMs.

- 3) Section 4 evaluates representative reasoning LLMs. Specifically, Section 4.1 reviews current mainstream reasoning benchmarks, covering both plain text and multimodal benchmarks across various task types. Section 4.2 outlines the current evaluation metrics, while Section 4.3 analyzes and compares the performance of mainstream reasoning LLMs with their foundational counterparts based on these benchmarks.
- 4) Section 5 highlights the limitations of existing reasoning LLMs and outlines several promising future development directions for these models.
- 5) Finally, we conclude the paper in Section 6 and provide a real-time tracking GitHub Repository to monitor the latest developments in the field.

We hope this survey serves as a valuable resource, fostering innovation and progress in this rapidly evolving domain.

1.2 Contribution of the Survey

Recently, several analyses and replications of specific technical approaches have been conducted [48]–[55], yet there remains a lack of systematic analysis and organization. Research [56] has focused only on slow-thinking methods during testing. Meanwhile, studies [57]–[59] have primarily concentrated on training or achieving reasoning LLMs, often from the perspective of RL.

Our survey distinguishes itself from and contributes to the existing literature in the following ways:

- 1) Rather than focusing on a single technical approach, we offer a comprehensive overview of the key concepts, methods, and challenges involved in reasoning LLMs.
- 2) We summarize the key advancements of early *System 2* and how they have paved the way for reasoning LLMs,

specifically in combination with foundational LLMs—a crucial aspect often overlooked in previous works.

- 3) We present a more thorough and inclusive summary of the core methods necessary for constructing reasoning LLMs, including but not limited to RL.

2 FOUNDATIONS OF REASONING LLMs

In this section, we provide a concise overview of the progress in foundational LLMs and the early development of key *System 2* technologies, highlighting critical advancements that, when combined with foundational LLMs, have paved the way for reasoning LLMs. These advancements include symbolic logic systems, MCTS, and RL.

2.1 Foundational LLMs

The development of foundational LLMs saw significant advancements with the introduction of pretrained Transformers [18] in 2018-2019, notably through BERT [19] and GPT [21]. These models leveraged unsupervised pretraining on vast text corpora, followed by fine-tuning for task-specific applications. This approach enabled them to develop a broad language understanding before specializing in tasks such as sentiment analysis, entity recognition, and question answering. BERT’s bidirectional context processing improved word understanding, while GPT excelled in text generation with its unidirectional design.

The release of GPT-2 [22] in 2019, with 1.5 billion parameters, marked a significant leap in generative performance, though it also raised ethical concerns. GPT-3 [23], with 175 billion parameters, further demonstrated the power of unsupervised pretraining, excelling in few-shot learning and performing well across a wide range of NLP tasks. In subsequent years, multimodal models like CLIP [60] and DALL-E [61] emerged, integrating text and visual inputs. These models enabled new tasks, such as generating images from text, and enhanced human-computer interaction.

By 2023-2024, models such as GPT-4 [62], LLaMA [25], and LLaVA [27] demonstrated advanced capabilities in reasoning, contextual understanding, and multimodal reasoning, processing both text and images. The evolution of foundational LLMs has revolutionized AI, enabling more sophisticated applications in language comprehension, problem-solving, and human-machine collaboration.

Summary: The development of foundational LLMs has progressed from pretrained transformers like BERT to multimodal models such as GPT-4, enhancing language understanding, text generation, and image processing. This advancement has led to significant breakthroughs in AI, improving language comprehension, problem-solving, and human-computer interaction. Building on deep learning advancements [63]–[66], foundational LLMs can learn extensive world knowledge and semantic relationships from vast textual or multimodal data. This enables them to exhibit emergent capabilities such as In-Context Learning (ICL) [67], prompt engineering [68], and Chain-of-Thought (CoT) reasoning [2], significantly enhancing their adaptability and creative problem-solving abilities.

Despite this progress, foundational LLMs operate similarly to *System 1* reasoning, relying on fast, heuristic-driven decision-making and lacking the step-by-step analysis characteristic of *System 2*. However, their developments lay a solid foundation for future reasoning LLMs—especially when integrated with the following early *System 2* technologies. This combination paves the way for more versatile, flexible, and human-like reasoning models.

2.2 Symbolic Logic Systems

Symbolic logic systems mark the earliest phase of AI, utilizing rules and logical principles to represent knowledge and draw conclusions [69], [70]. They are particularly effective in structured domains, where formal logic ensures precision.

Prolog, a logic programming language based on first-order logic, allows users to define facts, rules, and reason through queries. It has been pivotal in symbolic reasoning systems, especially in NLP and expert systems [71]–[73]. Logic-based systems like Prolog employ propositional and predicate logic for formal reasoning [74], [75]. From the 1960s to the early 1980s, this approach dominated AI, with systems like IBM’s LISP [76] for symbolic computation and Resolution Theorem Provers [77] for automated reasoning. In the 1970s, Marvin Minsky introduced Frames, which organized knowledge into structured frameworks, influencing both expert systems and cognitive science [78].

Summary: Symbolic logic systems were pivotal milestones in early AI development. Based on formal logic, they excelled in well-defined problems, particularly in structured environments. However, they also exposed the limitations of rigid, rule-based systems. Despite these constraints, symbolic logic remains foundational to the progress of AI.

Recent advancements in reasoning LLMs have greatly enhanced the emulation of human-like *System 2* cognitive processes through sophisticated thought architectures, known as Macro Action frameworks (Section 3.2.4). By combining symbolic templates or rules with foundational LLMs, macro actions have significantly improved their reasoning capabilities. Integrating macro actions into foundational LLMs has transformed their ability to handle complex

reasoning tasks, as hierarchical planning allows models to make high-level decisions before delving into specific problem details, mirroring symbolic logic’s structured approach.

2.3 Monte Carlo Tree Search

MCTS is a simulation-based search algorithm for decision-making and planning [79]. It constructs a search tree through four steps: *Selection*, which chooses the child node with the highest priority using the UCB1 formula:

$$UCB1 = \frac{w_i}{n_i} + c\sqrt{\frac{\ln N}{n_i}}, \quad (1)$$

where w_i is the total reward of node i , n_i is its visit count, N is the parent node’s visit count, and c balances exploration and exploitation. *Expansion* adds new nodes, *Simulation* performs random rollouts to evaluate them, and *Backpropagation* updates node statistics. MCTS has been widely used in tasks such as optimizing strategies in board games like Go [80] and in robotic path planning, where it helps robots navigate dynamic environments effectively [81].

Summary: MCTS has played a crucial role in the development of reasoning LLMs, particularly in Structural Search (Section 3.2.1). By simulating potential future reasoning paths and backpropagating estimated rewards, MCTS helps foundational LLMs efficiently identify the most promising, high-reward paths. This process mirrors human-like planning, where future consequences of decisions are considered before taking action. By dynamically exploring multiple reasoning trajectories, MCTS enables models to avoid getting stuck in suboptimal paths, making it easier to navigate complex decision spaces. This integration has significantly enhanced the ability of LLMs to handle intricate and dynamic reasoning problems, such as those requiring long-term planning or multi-step logical inferences. It has allowed LLMs to make more strategic and informed decisions, improving their overall performance in tasks that involve nuanced reasoning and strategic exploration.

2.4 Reinforcement Learning

RL is a type of machine learning where an agent learns to make decisions by interacting with an environment and receiving feedback in the form of rewards, aiming to maximize cumulative rewards over time [82]. Early breakthroughs in RL, such as Q-learning [83] and DQNs [84], revolutionized the field by enabling the handling of complex state spaces using Deep Neural Networks (DNNs) [85]. These methods paved the way for scaling RL to real-world tasks, where traditional tabular approaches fell short. The advent of deep RL marked a significant step forward, combining the power of deep learning with RL to process high-dimensional inputs, such as images and unstructured data.

A landmark achievement in deep RL was AlphaGo, which demonstrated RL’s potential by defeating a world champion in the complex game of Go through self-play [86]. This success highlighted deep RL’s ability to thrive in environments with large, continuous action spaces and uncertainty. Building on this, AlphaZero advanced the approach by mastering multiple board games—chess, Go, and Shogi—using self-play, MCTS, and DNNs [87]. AlphaZero’s ability to learn entirely from scratch, without prior human

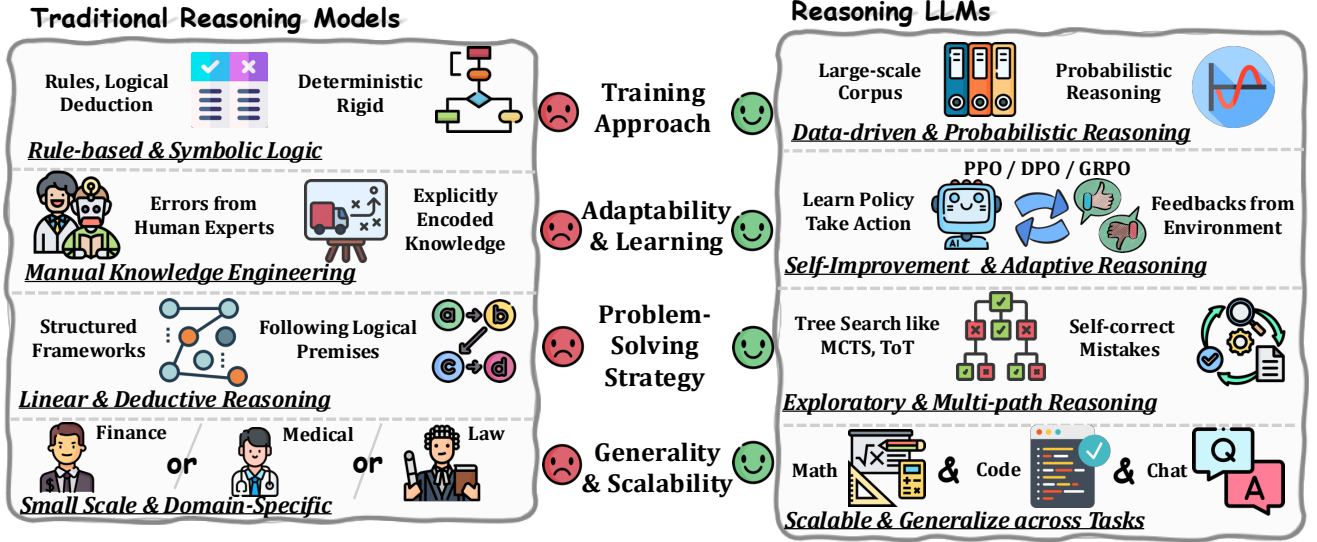


Fig. 3. A comprehensive comparison of traditional reasoning models and reasoning LLMs. Reasoning LLMs offer significant advantages over traditional models in areas such as training approaches, adaptability and learning, problem-solving strategies, and generality and scalability.

knowledge, showcased RL’s power in environments requiring long-term strategy and planning.

AlphaStar further expanded the boundaries of deep RL by excelling in the real-time strategy game StarCraft II. Unlike board games, StarCraft II presents dynamic, partially observable environments and demands multi-step, real-time decision-making [88]. AlphaStar’s success in this domain demonstrated deep RL’s capacity to adapt to complex decision-making scenarios that require both strategic planning and tactical execution. These advancements in RL and deep RL have greatly expanded AI’s potential, transitioning from well-defined, static environments to dynamic, complex settings that demand continuous learning and adaptation.

Summary: Deep RL has proven highly effective in solving complex decision-making tasks. AlphaGo exemplifies this by learning strategies through self-play and defeating the world champion in Go. This self-play concept laid the foundation for Self Improvement technology (Section 3.2.3) in reasoning LLMs, both relying on continuous feedback and adjustments to optimize strategies.

In RL, reward shaping has been crucial, especially for multi-step reasoning tasks [89]. By adjusting the reward signal to provide more granular feedback during intermediate steps, it helps agents navigate complex decision-making paths. This concept inspired the development of Reward Modeling (Section 3.2.2), particularly the process reward model, in reasoning LLMs. This model offers step-by-step supervision to identify and correct errors in the reasoning process. By mimicking human reasoning, the process reward model ensures more robust and interpretable results, especially in tasks like mathematical problem-solving and code generation, where step-by-step evaluation is critical.

Moreover, RL itself is a powerful tool for reasoning LLMs (Section 3.2.5). With a reward mechanism, RL guides foundational LLMs to find optimal solutions, especially in dynamic reasoning problems. Its simplicity and efficiency make RL invaluable for training and optimizing reasoning LLMs, enhancing the intelligence and self-evolution of AI

models. The integration of RL has led to significant advancements in reasoning LLMs, as demonstrated by DeepSeek-R1 [31], offering more flexible and efficient solutions.

3 BLUEPRINTING REASONING LLMs

In this section, we first analyze the features of reasoning LLMs from both output behavior and training dynamics perspectives. We then provide a detailed overview of the core methods that enable their advanced reasoning capabilities. Finally, we summarize the evolution of reasoning LLMs. A comprehensive comparison of traditional reasoning models and reasoning LLMs is shown in Figure 3.

3.1 Analysis of the Features of Reasoning LLMs

3.1.1 Output Behaviour Perspective

Explore and Planning Structure: Recent empirical studies have revealed that reasoning LLMs demonstrate a strong tendency for exploratory behavior in their output structures, especially when compared to models such as WizardMath [90] and DeepSeekMath [91], which primarily rely on conventional CoT reasoning approaches. This exploratory behavior is evident in the formulation of novel hypotheses and the pursuit of alternative solution paths. Research by [49] suggests that slow-thinking models engage in a latent generative process, particularly noticeable during the prediction of subsequent tokens. This claim is supported by [31], which observes that similar behaviors naturally arise during RL scale training. Furthermore, the Quiet-STaR framework [92] introduces an auxiliary pre-training phase focused on next-token prediction, highlighting the critical role of internal deliberation and exploratory mechanisms prior to content generation. Collectively, these findings underscore the complex and dynamic nature of reasoning processes in advanced LLMs, emphasizing the interaction between exploration and structured reasoning within their operational frameworks.

Verification and Check Structure: Analysis of OpenAI’s o1 [29] and o3 [30] models indicates that their reasoning frameworks incorporate both macro-level actions for long-term strategic planning and micro-level actions, including “Wait”, “Hold on”, “Alternatively”, and “Let’s pause”. These micro actions facilitate meticulous verification and iterative checking processes, ensuring precision in task execution. Such a dual-layered approach underscores the models’ capacity to balance overarching goals with granular, detail-oriented operations, thereby enhancing their overall functionality and reliability. To emulate this characteristic, Marco-o1 [93], during the MCTS process for constructing Long-CoT, assigns each tree node the state of “Wait! Maybe I made some mistakes! I need to rethink from scratch”, thereby facilitating the reflective nature of Long-CoT. Huatuo-o1 [94] employs a multi-agent framework to address the issue of incorrect CoT generation during validation. This is achieved by incorporating a prompt with “Backtracking” and “Correction” functionalities, which enables the correction process.

Longer Inference Length & Time: Recent research [49]–[52] indicates that reasoning LLMs often generate outputs exceeding 2000 tokens to tackle complex problems in coding and mathematics. However, this extended output length can sometimes lead to overthinking, where the model spends excessive time on a problem without necessarily improving the solution. Studies [49] highlight that while autoregressive generation and Classic CoT can effectively solve simpler problems, they struggle with more complex tasks. Research [95], [96] shows that in multimodal domains, many problems demand careful observation, comparison, and deliberation. Additionally, Search-o1 [97] suggests that slow-thinking mechanisms are particularly beneficial in areas requiring external knowledge or where potential knowledge conflicts arise. In medical scenarios, complex problems, such as those requiring test-time scaling techniques, demonstrate significant improvements [52].

Overly Cautious & Simple Problem Trap: Currently, reasoning LLMs have demonstrated strong performance in domains such as competitive-level mathematics [31], [54], [98], [99], complex coding [100], medical question answering [52], [94], and multilingual translation [93], [101]. These scenarios require the model to perform fine-grained analysis of the problem and execute careful logical reasoning based on the given conditions. Interestingly, even for straightforward problems like “ $2+3=?$ ”, reasoning LLMs can exhibit overconfidence or uncertainty. Recent research [102] notes that o1-like models tend to generate multiple solution rounds for easier math problems, often exploring unnecessary paths. This behavior contrasts with the lack of diverse exploratory actions for simpler questions, indicating a potential inefficiency in the model’s reasoning process.

3.1.2 Training Dynamic Perspective

Amazing Data Efficiency: Unlike traditional approaches that focus on expanding instruction sets with uniformly distributed difficulty levels, Studies [52], [54] suggest that constructing Slow-thinking CoT datasets with a focus on hard samples leads to better generalization in fields like medicine and mathematics. This approach diverges from the conventional practice of collecting diverse and evenly distributed instruction datasets.

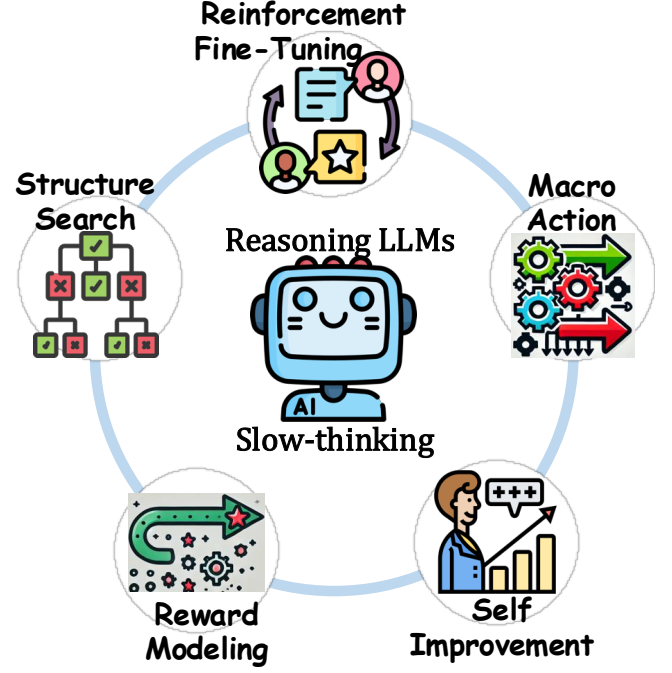


Fig. 4. The core methods enabling reasoning LLMs.

Sparse Training Method: Contrary to conventional wisdom, the development of effective reasoning LLMs does not require extensive datasets or dense reward signals. For example, STILL2 [51] demonstrated impressive performance using only 5,000 distilled samples, while Sky-T1 [99] achieved performance parity with QwQ [98] using just 17,000 Long-CoT samples. Similarly, RedStar [54] achieved exceptional results across both textual and multimodal tasks with only 4,000 core LongCoT samples. In comparison to simple CoT, Slow-thinking Supervised Fine-Tuning (SFT) data exhibits remarkable sample efficiency, often delivering comparable results with just 1/100th of the sample size. Additionally, research [103] emphasizes the significant training potential of online RL scaling algorithms, suggesting that non-dense RL supervision and even rule-based reward structures are sufficient for achieving high performance.

Parameter Characteristic: Training LLMs for slow-thinking, as characterized by the LongCoT approach, results in relatively uniform gradient norms across different layers. In contrast, fast-thinking, exemplified by the simplified CoT method, generates larger gradient magnitudes in the earlier layers, along with significant variability in gradient norms across layers. Empirical evidence suggests that larger models, particularly those exceeding 30 billion parameters, are more compatible with reasoning LLMs training due to their enhanced capacity for complex reasoning. Additionally, experiments conducted by RedStar [54] show that the benefits of data scaling vary across model sizes, with scaling effects being more pronounced and effective in larger models. This finding is supported by Deepseek-R1’s research [31], which demonstrates that a 670-billion-parameter model achieves performance metrics closely approximating those of the o1 benchmark, highlighting the scalability advantages of larger architectures in advanced reasoning tasks.

TABLE 1
Summary of Structure Search method.

	Category	Reasoning LLMs	Characteristic
Actions	Reasoning Steps as Nodes	RAP [14], ORM [104], Forest-of-Thought [105]	Actions represent intermediate reasoning steps.
	Token-level Decisions	CodeTree [106], SPaR [107], TreeBoN [108]	Actions involve generating tokens.
	Task-specific Structures	CWM [109], LLM-MCTS [110]	Actions are domain-specific.
	Correction and Exploration	RethinkMCTS [111], MCTSr [112]	Actions emphasize revisiting, refining, or backtracking to improve previous reasoning steps.
Rewards	Outcome-based Rewards	MC-NEST [113]	Correctness or validity of the final outcome.
	Stepwise Evaluations	RAP [14], SRA-MCTS [114]	Rewards are assigned at intermediate steps.
	Self-evaluation Mechanisms	SPaR [107], TreeBoN [108], MindStar [115]	Rewards rely on the model's own confidence.
	Domain-specific Criteria	LLM-MCTS [110], SR-MCTS [116]	Rewards are tailored to specific tasks.
	Iterative Preference Learning	LLaMA-Berry [117], Marco-o1 [93], ReST-MCTS* [118]	Rewards derive from comparing multiple solutions.

3.2 Core Method

In this section, we provide an overview of the core methods that drive the advanced reasoning capabilities of reasoning LLMs, as shown in Figure 4. These include Structure Search, Reward Modeling, Self Improvement, Macro Action, and Reinforcement Fine-Tuning. We also highlight representative reasoning LLMs for each method.

3.2.1 Structure Search

Reasoning LLMs aim to achieve high accuracy and depth in solving complex problems by emulating the deliberate and methodical nature of human reasoning. However, despite recent advancements, current foundational LLMs face inherent limitations when addressing intricate reasoning tasks. These limitations arise from their lack of an internal world model to simulate environmental states, their inability to predict the long-term outcomes of reasoning paths, and their failure to iteratively refine reasoning steps based on future states or rewards [8]. As a result, these shortcomings hinder foundational LLMs from effectively balancing exploration and exploitation in vast reasoning spaces, creating challenges in tasks that require multi-step reasoning, such as complex mathematics, logical inference, or strategic decision-making [119].

MCTS, a powerful search and optimization algorithm, effectively addresses these challenges by providing a structured framework to explore and evaluate reasoning paths systematically. It operates by constructing a reasoning tree, where each node represents a reasoning state, and actions expand the tree by considering potential next steps. Through the simulation of future states and the iterative backpropagation of estimated rewards, MCTS allows foundational LLMs to efficiently identify high-reward reasoning paths, mirroring human planning processes. This approach aligns with the core principles of reasoning LLMs, where thorough analysis and deliberate exploration are essential for generating well-reasoned outputs. Recent methods, such as RAP [14], enhance foundational LLMs by integrating MCTS with a world model, enabling the system to iteratively refine intermediate reasoning steps and improve future predictions. Similarly, Forest-of-Thought [105] utilizes MCTS to dynamically explore multiple reasoning trajectories, revisiting flawed paths and refining outcomes.

The application of MCTS in reasoning tasks extends beyond traditional problem-solving to highly specialized domains. For example, frameworks like SRA-MCTS [114] and MC-NEST [120] showcase the utility of MCTS in tackling technical challenges such as code generation and mathematical reasoning, where intermediate steps are iteratively

evaluated and refined. In fields like instructional alignment, frameworks such as SPaR [107] and Marco-o1 [93] leverage MCTS to refine responses and align reasoning trajectories with human preferences or desired outcomes. Additionally, task-specific implementations like HuatuoGPT-o1 [94] underscore MCTS's crucial role in navigating highly specialized domains, such as medical reasoning, where accuracy and robustness are paramount.

MCTS also enables models to go beyond single-pass reasoning methods, such as CoT or Tree-of-Thought, by incorporating mechanisms to revisit, critique, and refine reasoning steps dynamically [111], [121]. This iterative capability is essential for tackling tasks with vast decision spaces or those requiring long-term planning, where earlier decisions can significantly impact final outcomes. By allowing LLMs to simulate, evaluate, and refine multiple reasoning paths, MCTS introduces a level of adaptability and strategic exploration that traditional approaches lack. As shown by AlphaZero-like tree-search [104] and Search-o1 [97], MCTS enables reasoning LLMs to not only achieve better performance on specific tasks but also exhibit enhanced generalization capabilities across diverse domains.

The integration of MCTS into LLMs depends on defining actions and rewards to guide reasoning path exploration and assess quality. As shown in Table 1, we classify the actions in prior work into four categories:

- 1) **Reasoning Steps as Nodes:** Actions represent intermediate reasoning steps or decisions, such as selecting rules, applying transformations, or generating sub-questions [14], [104], [105], [119].
- 2) **Token-level Decisions:** Actions involve generating tokens or sequences (e.g., the next word, phrase, or code snippet) [106]–[108], [122].
- 3) **Task-specific Structures:** Actions are domain-specific, such as moving blocks in blocksworld, constructing geometry in geometry problem-solving, or modifying workflows in task planning [109], [110], [123].
- 4) **Self-correction and Exploration:** Actions focus on revisiting, refining, or backtracking to improve previous reasoning steps [111], [112], [124].

Additionally, as shown in Table 1, we classify the reward design into five categories:

- 1) **Outcome-based Rewards:** Rewards focus on the correctness or validity of the final outcome or solution, including the validation of reasoning paths or task success [113], [119], [123].
- 2) **Stepwise Evaluations:** Rewards are assigned at intermediate steps based on the quality of each step or its contribution toward the final outcome [14], [104], [114].

TABLE 2
Summary of Reward Modeling method.

Category	Methods	Data Source	Model Refinement		Applications	Characteristic
			Strategy	Learning		
ORM	DIVERSE [127]	Prompting	Fine-tuning	SFT	Multiple Reasoning Tasks	Weighted Voting Verifier
	MATH-SHEPHERD [128]	Sampling	Feedback-guided	SFT & RL	Math Reasoning	Correctness Score Assignment
	AutoPSV [129]	Prompting	Feedback-guided	SFT	Math / Commonsense Reasoning	Automated Process Supervision
	Implicit PRMs [130]	Sampling	Fine-tuning	SFT & RL	Math Reasoning	Obtaining PRM from ORM
MCTS	OVM [131]	Sampling	Feedback-guided	SFT	Math Reasoning	Guided Decoding
	ReST-MCTS* [132]	Sampling	Self-training	SFT & RL	Multiple Reasoning Tasks	MCTS and Self-training
	OmegaPRM [133]	MCTS with Binary Search	Feedback-guided	SFT	Math Reasoning	Divide-and-Conquer MCTS
	ReARTeR [134]	Sampling	Feedback-guided	SFT & RL	QA	Retrieval-Augmented Generation
	Consensus Filtering [135]	MCTS Data Construction	Feedback-guided	SFT	Math Reasoning	Consensus Filtering Mechanism
PRM	ORPS [136]	Sampling	Feedback-guided	SFT	Code Generation	Supervising Outcome Refinement
	Step-DPO [137]	Sampling	Feedback-guided	SFT & RL	Math Reasoning	Step-wise Preference Pairs
	AdaptiveStep [138]	Response Dividing	Feedback-guided	SFT	Math Reasoning, Code Generation	Dividing Reasoning Steps

- 3) **Self-evaluation Mechanisms:** Rewards rely on the model's own confidence or self-assessment (e.g., likelihood, next-word probability, or confidence scores) [107], [108], [115].
- 4) **Domain-specific Criteria:** Rewards are tailored to specific tasks, such as symmetry and complexity in geometry or alignment with human preferences in text generation [110], [116], [123].
- 5) **Iterative Preference Learning:** Rewards are derived from comparing multiple solutions or reasoning paths, guiding learning dynamically [93], [117], [118].

Summary: Despite its advantages, structure search-based (i.e., MCTS) reasoning LLMs often suffer from substantial computational overhead due to the large number of simulations required. This makes them less suitable for tasks that demand real-time decision-making or operate under resource constraints [125]. Additionally, the effectiveness of MCTS is highly dependent on well-designed reward mechanisms and action definitions, which can vary significantly across different domains, thus posing challenges to its generalizability [126].

3.2.2 Reward Modeling

Two primary training paradigms are used to tackle multi-step reasoning tasks: outcome supervision and process supervision. Outcome supervision emphasizes the correctness of the final answer at a higher level of granularity, and the resulting model is referred to as the Outcome Reward Model (ORM) [32], [139]. In contrast, process supervision provides step-by-step labels for the solution trajectory, evaluating the quality of each reasoning step. The resulting model is known as the Process Reward Model (PRM) [37], [140], [141]. The main distinction between ORM and PRM is illustrated in Figure 5.

PRM offers significant advantages [128], [142] in complex reasoning tasks for several key reasons. First, it provides fine-grained, step-wise supervision, allowing for the identification of specific errors within a solution path. This feature is especially valuable for RL and automated error correction. Second, PRM closely mirrors human reasoning behavior, which relies on accurate intermediate steps to reach correct conclusions. Unlike ORM, PRM avoids situations where incorrect reasoning can still lead to a correct final answer, thus ensuring more robust and interpretable reasoning. While PRM has primarily been applied to complex

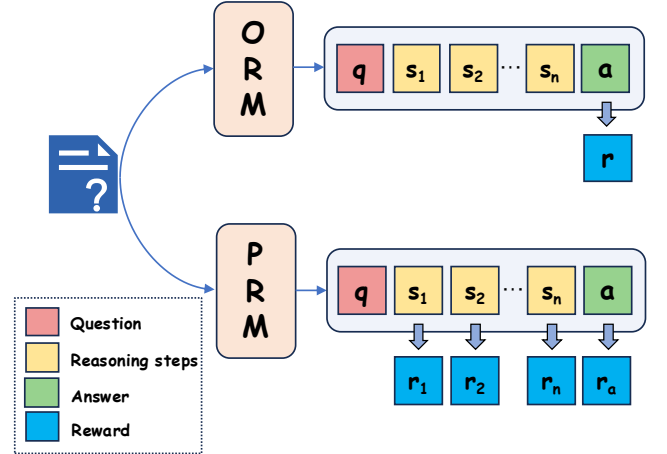


Fig. 5. The comparison between ORM and PRM for assessing a complete solution trajectory. ORM only provides a single reward based on the correctness of the final answer, while PRM evaluates the quality of each reasoning step throughout the process.

mathematical problems, its benefits have recently driven applications in other fields. For instance, ORPS [136] utilizes PRM to address complex code generation challenges, while Step-DPO [137] combines process supervision with the Direct Preference Optimization (DPO) algorithm [143] to improve long-chain mathematical reasoning. A summary of Reward Modeling method is presented in Table 2.

Summary: Despite the advantages of PRMs, they present several challenges. The primary difficulty is obtaining process supervision-labeled data, which is often both costly and time-consuming. To address concerns related to scalability, efficiency, and accuracy, researchers have explored various automated annotation methods. For example, MATH-SHEPHERD [128] utilizes the correctness of the final answer to define the quality of intermediate steps based on their potential to lead to the correct outcome, automating the step-wise data collection process. ReST-MCTS* [132] combines process reward guidance with MCTS to generate higher-quality reasoning traces through extensive rollouts. Similarly, OmegaPRM [133] employs the MCTS framework while introducing a divide-and-conquer algorithm for automated process supervision data generation. Another novel approach involves using ORM to train a PRM. Yuan et al. [130] propose training a PRM implicitly by leveraging

TABLE 3
Summary of Self Improvement method.

Stage	Methods	Data Source	Model Refinement		Applications
			Feedback	Strategy	
Training	STaR [144]	Few-shot	Language Model	SFT	QA, Arithmetic Reasoning
	Quiet-STaR [92]	Token-level Exploration	Language Model	RL	QA, Arithmetic Reasoning
	V-STaR [145]	Sampling	Verifier	SFT	Arithmetic Reasoning, Code Generation
	B-STaR [146]	Sampling	Reward Model	SFT	Arithmetic Reasoning, Code Generation
	rStar-Math [147]	MCTS Data Construction	Reward Model	SFT	Arithmetic Reasoning
	ReST [148]	Sampling	Reward Model	RL	Machine Translation
	ReST-EM [149]	Sampling	Language Model	EM for RL	Arithmetic Reasoning, Code Generation
	ReST-MCTS* [132]	Sampling	Reward Model	SFT, RL	Reasoning
	ENVISIONS [150]	Sampling	Environment Guided	SFT	Web Agents, Reasoning
	RISE [151]	Sampling	Reward Function	Weighted SFT	Arithmetic Reasoning
	STIC [152]	Few-shot	Language Model	SFT	Vision Language Model Tasks
	SIRLC [153]	Question Answering	Language Model	RL	Reasoning, Translation, Summary
	AlpacaFarm [154]	Existing Data	Language Model	SFT	None (Intrinsic Evaluation)
Inference	Self-Refine [155]	Independent of Training Data	Language Model	Few-shot Demonstration	Code Generation, Sentiment Reversal, Acronym Generation
	Self-Check [156]	Independent of Training Data	Language Model	Step Check	QA, Arithmetic Reasoning
	CRITIC [157]	Independent of Training Data	Language Model	External Tools	QA, Arithmetic Reasoning, Detoxification
	ROSE [158]	Independent of Training Data	Language Model	Distributed Prompt	Safety, Knowledge
	Self-Verification [159]	Independent of Training Data	Language Model	Re-Ranking	Arithmetic Reasoning
	SelfEval-Decoding [160]	Independent of Training Data	Language Model	Beam Search	Arithmetic/Symbolic Reasoning
	IPS [161]	Independent of Training Data	Language Model	Constrained Decoding	Dialogue
	Control-DAG [162]	Independent of Training Data	Language Model	Constrained Decoding	Dialogue, Open-domain Generation
	Look-Back [163]	Independent of Training Data	Language Model	Contrastive Decoding	Alleviating Repetitions
	LeCo [164]	Independent of Training Data	Language Model	Constrained Decoding	QA, Reasoning

ORM training on cheaper datasets, under mild reward parameterization assumptions. They also provide theoretical guarantees for the performance of this implicit PRM, demonstrating its practicality and cost-effectiveness.

In addition to data collection, PRMs face challenges related to trustworthiness [134], categorized as follows:

- 1) **Lack of Explanations:** Current PRMs often generate scores for reasoning steps without sufficient explanations, limiting interpretability and hindering their usefulness in refining reasoning during test-time.
- 2) **Bias in Training Data:** Data collection methods, such as MCTS, tend to introduce distributional biases, assigning disproportionately higher scores to the majority of questions. As a result, PRMs struggle to effectively identify erroneous reasoning steps.
- 3) **Early-Step Bias:** PRMs show lower accuracy in predicting rewards for earlier reasoning steps compared to those closer to the final answer. This issue stems from the increased randomness and uncertainty associated with the initial steps in the reasoning process.

3.2.3 Self Improvement

Reasoning LLMs exemplify a progression from weak to strong supervision, while traditional CoT fine-tuning faces challenges in scaling effectively. Self improvement, using the model’s exploration capabilities for self-supervision, gradually enhances LLMs performance in tasks such as translation [148], mathematical reasoning [144], [149], and multimodal perception [152]. This approach fosters exploration and application within reasoning LLMs [147], [165]. A summary of Self Improvement method is presented in Table 3.

Training-based self improvement in LLMs can be categorized based on exploration and improvement strategies. The exploration phase focuses on data collection to facilitate subsequent training improvements, with notable variations in approach. STaR [144] uses few-shot examples for data gathering, while ReST [148], ReST-EM [149], and ENVISIONS [150] rely on multiple samplings of complete

trajectories. Quiet-STaR [92] explores at the token level, introducing concepts like meta-tokens and non-myopic loss to enhance supervision. Additionally, ReST-MCTS* [132] and rStar-Math [147] generate training data through MCTS.

Improvement strategies also exhibit significant diversity. For instance, STaR and its derivatives, such as V-STaR [?] and B-STaR [146], combine filtering with SFT. ReST and its variants typically introduce innovative reward calculation methods to enhance RL training for policy models. RISE [151] incorporates external feedback, recording rewards and refining responses through distillation during the improvement process. Notably, rStar-Math [147] demonstrates that small models have achieved *System 2* reflective capabilities through self-evolving training approaches.

Test-time self improvement leverages the consistency of a model’s internal knowledge to correct hallucinations during inference. These approaches can be categorized into three main types: methods that refine answers using prompts [155], [156], approaches that utilize external tools [157], and techniques that leverage logits without the need for external tools or prompts [163], [164].

3.2.4 Macro Action

Recent advancements in LLMs have driven progress in emulating human-like *System 2* cognitive processes via sophisticated thought architectures, often referred to as macro action frameworks. These structured reasoning systems go beyond traditional token-level autoregressive generation by introducing hierarchical cognitive phases, such as strategic planning, introspective verification, and iterative refinement. This approach not only enhances the depth of reasoning but also broadens the solution space, enabling more robust and diverse problem-solving pathways. A summary of Macro Action method is presented in Table 4.

We classify the progress of macro action into two aspects:

- 1) **Test-time Scaling through Macro Action Operationalization:** Recent research identifies two key methodologies for improving reasoning performance during

TABLE 4
Summary of Macro Action method.

Methods	Usage	Action Attribute			Reflection	Modality	Representative Action
		Action Source	Action Number	Learning			
Self-Check [166]	Verification	Human-Designed	4	ICL	✓	✓	Target Extraction, Information Collection, Step Regeneration, Result Comparison
LeMa [167]	Synthetic Data	Human-Designed	3	ICL & SFT	✓	✓	Incorrect Step Recognition, Explanation, Correct Solution:
REFINER [168]	Verification/Exploration	Human-Designed	2	ICL & SFT	✓	✓	Critic, Generate
HiICL-MCTS [169]	Exploration	Human-Designed	5	ICL	✓	✓	System Analysis, One-Step Thought, Divide and Conquer, ..., Self-Reflection and Refinement
SUPERCORRECT [170]	Distill	In-Context Learning	Dynamic	SFT & RL	✓	✓	—
ReasonFlux [171]	Synthetic Data/Exploration	Human-Designed	~500	ICL & SFT & RL	✓	✓	—
rStar [172]	Exploration	Human-Designed	5	ICL & RL	✓	✓	One-step thought, Propose Next Sub-question & Answer, ..., Rephrase question
LLaMA-Berry [173]	Exploration	Human-Designed	2	ICL & RL	✓	✓	Reflection, Error Re-correction
Huatuo-o1 [94]	Synthetic Data	Human-Designed	4	ICL & SFT	✓	✓	Backtracking, Exploring New Paths, Verification, Correction
Marco-o1 [93]	Verification	Human-Designed	1	ICL & SFT	✓	✓	Reflection
BoT [174]	Exploration	In-Context Learning	Dynamic	ICL	✓	✓	Solving Quadratic Equation, Array Sorting, ..., Search Algorithms)
rStar-Math [147]	Exploration	In-Context Learning	1	ICL & RL	✓	✓	Python comment
Mulberry [175]	Synthetic Data	In-Context Learning	1	ICL & SFT	✓	✓	Reflection
LLaVA-CoT [176]	Synthetic Data/Exploration	Human-Designed	4	SFT	✓	✓	Summary, Caption, Reasoning, Conclusion
LLaMAV-o1 [177]	Verification/Exploration	Human-Designed	4173	Curriculum Learning	✓	✓	Detailed Caption Generation, Logical Reasoning, ... Final Answer Generation
AtomThink [178]	Synthetic Data/Exploration	In-Context Learning	>100	SFT & RL	✓	✓	Variable Definition, Calculations, Graphs Analysis, ..., Verification
RedStar [54]	Distill	Human-Designed	2	SFT	✓	✓	Wait, Alternately
Auto-CoT [179]	Exploration	In-Context Learning	2	ICL	✓	✓	Question clustering, Demonstration Sampling
PoT [180]	Verification	In-Context Learning	1	ICL	✓	✓	Code language conversion
PAL [181]	Verification	In-Context Learning	1	ICL	✓	✓	Code language conversion
Decomposed Prompt [182]	Exploration	Human-Designed	3	ICL	✓	✓	Problem Split, Subproblem Solving, Answer Merge
Least-to-Most [183]	Exploration	Human-Designed	2	ICL	✓	✓	Problem Decomposition, Subproblem Solving

inference and test-time scaling. HiICL-MCTS [169] employs a deliberate search through seed data to generate action-chain templates consisting of macro actions, thereby facilitating an action-chain-guided approach to test-time reasoning. ReasonFlux [171] utilizes an iterative test-time scaling framework, harnessing external high-level thought templates to iteratively refine and update the current CoT.

2) Macro Action-Enhanced Data Synthesis Paradigms: A key application of macro actions in complex reasoning is in the synthesis of reasoning data. In data synthesis and training frameworks, macro action architectures enhance reasoning diversity and generalization. Recent research has shown that integrating or synthesizing a CoT process with macro actions within the reasoning sequence can significantly improve the data efficiency of the reasoning chain. For instance, LLaVA-CoT [176] enhances CoT data synthesis by externalizing intermediate reasoning steps across multiple modalities. AtomThink [178] generates the AMATH-SFT dataset using a structured g1 prompt [184], achieving superior performance on long-horizon reasoning tasks compared to traditional CoT approaches. CoAct [185] introduces a dual-agent collaborative reasoning framework, where a global planning agent executes overarching macro-actions, while a local execution agent carries out specific sub-actions within those broader actions.

Macro actions also play a crucial role in enhancing Self-improvement frameworks. rStar-Math [147] utilizes high-level deliberate search through Code-augmented CoT, generating diverse and reliable solutions while achieving proactive search capabilities. Satori [186] integrates CoT with RL, incorporating “<reflect>”-style macro actions to diversify exploration and alleviate policy saturation in online RL environments. Huatuo-o1 [94] combines hierarchical planning with domain-specific knowledge bases to improve medical reasoning. Additionally, ReasonFlux [171] dynamically re-configures reasoning templates (e.g., breaking down calculus problems into symbolic and numeric phases) to align with the problem structure.

3.2.5 Reinforcement Fine-Tuning

Reinforcement Fine-Tuning (RFT) [187] is an innovative technique recently introduced by OpenAI, designed to enable developers and engineers to fine-tune existing models

for specific domains or complex tasks. Unlike general SFT, RFT focuses on optimizing the model’s reasoning process by using a reward mechanism to guide the model’s evolution, thereby enhancing its reasoning capabilities and accuracy. The core of RFT lies in improving the model’s performance in a specific domain with minimal high-quality training data [188], an appropriate reward model [189], and a stable optimization process [190]. A summary of RFT method is presented in Table 5.

DeepSeek-R1 [31], which employs a verifier reward-based strategy, has shown significant performance improvements compared to traditional methods like SoS [191]. Key advantages include:

- 1) Simplified Training Pipeline:** RL supervision streamlines data construction and training processes, eliminating the need for complex stepwise search mechanisms.
- 2) Enhanced Scalability:** Online RL training facilitates efficient scaling on large datasets, particularly for complex reasoning tasks.
- 3) Emergent Properties:** DeepSeek-R1 [31] demonstrates unique emergent capabilities, such as Long-CoT reasoning, which are difficult to achieve through SFT alone.

Despite its strengths, RFT faces the following challenges:

- 1) Unclear Mechanism behind Reasoning:** The underlying mechanisms driving the reasoning improvements in DeepSeek-R1 remain poorly understood. For example, while DeepSeek-R1 exhibits emergent properties (e.g., “Emergent Length Increasing”, “Aha moments”), studies such as [219] suggest that capabilities like Long-CoT might already exist in the base model, rather than solely emerging from RL training. Furthermore, performance gains observed in smaller models (e.g., Qwen-Math-2B/7B [220]) occur without noticeable “Aha moments”, complicating causal interpretations.
- 2) Reward Model Saturation:** Many existing RL algorithms face reward model saturation, typically manifested as exploration collapse after around 100 training steps. Although DeepSeek-R1 alleviates this issue through specialized reward formatting, methods like ReFT [189] and Satori [186] propose alternating sampling and SFT distillation to combat reward hacking and exploration collapse.
- 3) Unstable Long-CoT Generation:** Long reasoning chains generated by RFT are prone to instability, including context overflow, failure to return final answers, and

TABLE 5
Summary of RFT method.

Methods	Model Attribute		Modality	Reward Type	Incentivize Attribute			Application & Benchmark
	Foundational LLMs				Algorithm	Learning	Incentivize Sample	
Reason RFT Project								
DeepSeek-R1-Zero [31]	DeepSeek-V3		Rule-Outcome-Reward	GPRO	RL	800K	Multiple Tasks	
DeepSeek-R1 [31]	DeepSeek-V3		Rule-Outcome-Reward	GPRO	RL & SFT	800K	Multiple Tasks	
Kimi v1.5 [192]	–		Rule-Outcome-Reward	PPO*	RL & SFT	–	Multiple Tasks	
ReFT [189]	Galactica, CodeLLama		Rule-Outcome-Reward	PPO*	RL & SFT	3k/7k/8k/15k	GSM8k/SVAMP/MathQA	
RFT [193]	LLaMA-3-3/8B-Instruct,Qwen-2.5-7B-Instruct		Rule-Outcome-Reward	Reinforce++	RL & SFT	1.2K	Multiple Math Task	
Satori [186]	Qwen-2.5-Math-7B		Rule-Outcome-Reward	PPO	RL & SFT	66K	Multiple Math Task	
QCLASS [194]	Llama-2-7B-Chat		Process-Reward	QNet	RL & SFT	1.9K/1.5K/3.3K	WebShop, ALFWorld, SciWorld	
PRIME [195]	Qwen2.5-Math-7B		Rule-Process-Outcome-Reward	PPO	RL & SFT	150K	Math, Code Tasks	
DeepScaleR [196]	DeepSeek-R1-Distill-Qwen-1.5B		Rule-Outcome-Reward	Iteratively GPRO	RL	40K	Multiple Math Task	
PURE [197]	Qwen2.5-Math-7B		Rule-Process-Outcome-Reward	PPO+RLOO	RL	8K	Multiple Math Task	
SimpleRL [103]	Qwen2.5-Math-7B		Rule-Outcome-Reward	PPO	RL	8K	Multiple Math Task	
Open-R1 [198]	Qwen2.5-1.5B-Instruct		Rule-Outcome-Reward	GPRO	RL & SFT	8K	Multiple Math, Code Task	
TinyZero [199]	Qwen2.5-0.5B/3B		Rule-Outcome-Reward	GPRO	RL	–	CountDown Task	
Ota-Zero [200]	Qwen-2.5-Series, DeepSeek-Series, Rho, Llama-3.x		Rule-Outcome-Reward	GRPO	RL	0.5K	CountDown Task	
Ota [201]	RHO-1b/Qwen2.5-3B		Rule-Outcome-Reward	GPRO/PPO	RL	7.5K	GSM8K	
LIMR [202]	Qwen-Math-7B		Rule-Outcome-Reward	PPO	RL	1.3K	Multiple Math Task	
Critic-RL [203]	Qwen2.5-Coder-32B		Rule-Outcome-Reward	GPRO*	RL & SFT	18.8K	Multiple Code Task	
Logic-R1 [204]	Qwen2.5-7B-Instruct-1M		Rule-Outcome-Reward	REINFORCE++*	RL	5K	Multiple Math, Logic Task	
Online-DPO-R1 [205]	Qwen2.5-MATH-7B		Rule-Outcome-Reward	DPO	RL & SFT	207.5K	Multiple Math Task	
OpenReason-Zero [206]	Qwen-2.5-7B/32B		Rule-Outcome-Reward	PPO	RL	57K	Multiple Math Task, GPQA, MMLU	
RLHF-V [207]	OmniLM-12B		Process-Reward	DDPO	RL	1.4K	Multiple Tasks	
RLAIF [208]	PaLM 2 Extra-Small		Rule-Outcome-Reward	RLAIF	RL	–	Summary and Conversation Generation	
MM-RLHF [209]	LLaVA-onevision-7B		Process-Reward	MM-DPO	RL	120K	MM-RLHF-RewardBench/SafetyBench	
Align-DS-V [210]	LLaVA-v1.5-7B,Qwen2-VL		Process-Reward	PPO, DPO	RL & SFT	200K	Align-Anything, Eval-Anything	
R1V [211]	Qwen2-VL,Qwen2.5-VL		Rule-Outcome-Reward	GRPO	RL	70K/70K/8K	Multiple Tasks	
VLM-R1 [212]	Qwen2.5-VL		Rule-Outcome-Reward	GRPO	RL	120K	Multiple Tasks	
LMM-R1 [213]	Qwen2.5-VL		Rule-Outcome-Reward	PPO/RLOO	RL	8K	Multiple Tasks	
Open-R1-Video [214]	Qwen2-VL-7B		Rule-Outcome-Reward	GRPO	RL	4K	Multiple Tasks	
Easy-R1 [215]	Qwen2.5-VL		Rule-Outcome-Reward	GRPO	RL	3K	Multiple Tasks	
Analysis RFT Project								
Demystify-LongCoT [216]	Llama-3.1-8B, Qwen2.5 -7B-Math		Rule-Outcome-Reward	PPO/Reinforce++	RL & SFT	7.5K	Multiple Math, MMLU	
RLHF-Scale [217]	GLM4-9B		Process-Reward	PPO	RL	11K	Multiple Tasks	
MD-CoT [218]	–	–	Process-Reward	PPO	RL	–	–	

sensitivity to reward shaping [102]. For instance, methods like [216] inadvertently introduce cosine reward functions, which degrade performance with increased iterations. O1-Prune [221] uses post-hoc length pruning techniques [192] (via RL/SFT) to stabilize outputs.

Future directions for RFT may include several exciting and innovative advancements, such as:

- 1) **Efficient and Stable RL Frameworks:** There is a need to develop more robust RL algorithms that prevent reward saturation and exploration collapse. [216] reveals that REINFORCE++ [222] underperforms when combined with KL divergence regularization, suggesting the need for alternative methods. Future work should revisit classic RL algorithms in the context of modern LLMs training to optimize both stability and efficiency.
- 2) **Scaling RFT:** Current RL-Supervise models rely on curated, verifiable prompts selected from large-scale datasets. Future research should focus on synthesizing high-quality, diverse prompts to improve generalization. [217] shows that merely scaling policy/reward models or increasing sample sizes results in diminishing returns, while expanding the scope of PRM and R1 training data holds greater promise. Hybrid approaches, such as combining RL with SFT or curriculum learning, should be explored to enhance scalability.
- 3) **Controlling Long-CoT Stability:** Adaptive reward shaping mechanisms are needed to balance reasoning length, coherence, and answer correctness. Techniques such as O1-Prune [221] demonstrate the value of post-hoc length regularization, but dynamic in-training controls are necessary. Hierarchical RL frameworks should be investigated to decompose long reasoning chains into manageable sub-tasks, reducing instability.
- 4) **Theoretical and Empirical Analysis:** It is essential to clarify the relationship between RL training and the capabilities of the base model. For instance, it should be determined whether emergent properties (e.g., Long-

CoT) arise from RL optimization or are latent traits of the base model. Systematic studies on reward design principles (e.g., sparse vs. dense rewards, multi-objective balancing) should be conducted to avoid unintended behaviors such as reward hacking.

Summary: RFT presents a promising direction for advancing LLMs reasoning, as evidenced by DeepSeek-R1 [31]. However, challenges such as reward saturation, unstable long reasoning chains, and unclear emergent mechanisms require urgent attention. Future efforts should prioritize algorithmic innovation, scalable prompt synthesis, and theoretical grounding to fully unlock the potential of RL-driven reasoning LLMs.

3.3 Evolutionary of Reasoning LLMs

The evolution of reasoning LLMs has progressed by several distinct stages, with various strategies developed to overcome the limitations of direct autoregressive inference and build more advanced slow-thinking reasoning architectures.

In the early stages, reasoning LLMs primarily focused on enhancing pre-trained LLMs with external reasoning algorithms, without altering the underlying model parameters. Approaches such as Tree of Thoughts [223] and Reasoning via Planning [14] utilized LLMs-driven Breadth-First Search, Depth-First Search, and MCTS [79], [105], [108], [224] to simulate human-like reasoning processes. These methods represented reasoning as tree or graph traversals, where intermediate reasoning states were depicted as nodes, and various reasoning strategies produced distinct reasoning paths. The final decision was made through additional voting mechanisms [3] or Monte Carlo-based value estimation to identify the optimal path.

However, these externalized slow-reasoning approaches introduced several challenges:

- 1) **Limited Exploration Space:** The search-based methods required predefined constraints on the breadth, depth,

TABLE 6
Statistics of benchmarks for reasoning LLMs.

Domain	Benchmark	Venue	Language	Size	Level
Math	AIME 2024 [226]	-	English	30	Competition
	MATH-500 [37]	ICLR 2024	English	500	Competition
	AMC 2023 [227]	-	English	30	Competition
	Olympiad Bench [228]	ACL 2024	English/Chinese	8,476	Competition
Code	Codeforces	-	English	-	Expert
	SWE-bench [229]	ICLR 2024	English	2,294	Expert
	LiveCodeBench [230]	ArXiv 2024	English	-	Expert
Science	GPQA Diamond [231]	COLM 2024	English	448	University
	MMLU-Pro [232]	NeurIPS 2024	English	12,032	Hybrid
Agent	WebShop [233]	NeurIPS 2022	English	1,600	Hybrid
	WebArena [234]	ICLR 2024	English	812	Hybrid
	SciWorld [235]	EMNLP 2022	English	7,200	Hybrid
	TextCraft [236]	NAACL 2024	English	200	Hybrid
Medicine	JAMA Clinical Challenge [237]	NAACL 2025	English	1,524	Expert
	Medbullets [237]	NAACL 2025	English	308	Expert
	MedQA [238]	ArXiv 2020	English/Chinese	61,097	Expert
Multimodality	MMMU [239]	CVPR 2024	English	11,500	Hybrid
	MathVista [240]	ICLR 2024	English	6,141	Middle School
	MathVision [241]	NeurIPS 2024	English	3,040	Middle/High School
	CMMaTH [242]	COLING 2025	English/Chinese	23,856	Middle/High School
	PGPS9K [243]	IJCAI 2023	English	9,023	Middle School

and granularity of the search space, which often restricted the LLM’s exploration to a narrow reasoning space. Furthermore, the reasoning strategies across different child nodes of the same parent node frequently lacked sufficient diversity, further limiting exploration.

- 2) **Limited Experience Sharing:** Exploration experiences and reasoning information across different paths could only be assessed based on reward models or self-consistency among outcomes. Additionally, search-based methods significantly increased computational overhead, relying on reward models such as PRM/ORM for tree pruning or speculative decoding techniques to accelerate inference.

To overcome these limitations, subsequent models such as rSTaR [172], LLaMAV-o1 [177], HiICL-MCTS [169], Mulberry [175], g1 [184], and Thinking-Claude [225] introduced richer action spaces. These enhanced action spaces offered high-level planning cues, broadening the model’s exploration scope and enabling more comprehensive structured search processes. However, this approach necessitated careful design of the action spaces to ensure their effectiveness.

With the introduction of models like o1 [29] and QwQ [98], external reasoning paradigms were internalized within the LLM’s context. These models initially performed exploratory macro-planning to generate an initial reasoning path, followed by contextual exploration of alternative paths. Through mechanisms like “Rethink” and “Verification”, these models produced extended reasoning chains. To replicate this internalized capability, STILL-1 [224] linearized tree search outputs into long reasoning chains with attributes such as “Rethink”, “Wait”, and “Explore New Path”. Similarly, STILL-2 [53] and sky-T1 [99] synthesized long reasoning chains using distillation techniques. However, the linearized reasoning chains derived from search-based methods struggled to match the quality of those produced by distillation approaches.

Recent advancements, including DeepSeek-R1 [31] and

Kimi-k1.5 [192], have demonstrated the potential of RL to enhance models like DeepSeek-V3 [17], resulting in the emergence of complex behaviors such as long reasoning chains, reflective reasoning, and advanced planning capabilities. Remarkably, these sophisticated behaviors were achieved through simple RL scaling. SimpleRL [103] sought to replicate these capabilities using a streamlined pipeline and minimal codebase, while R1V [211] explored the development of multimodal reasoning models based on multimodal foundation architectures.

Summary: The evolution of reasoning LLMs has shifted from externally augmented reasoning to internally embedded reasoning. Recent developments emphasize the potential of RL-based scaling to unlock advanced capabilities.

4 BENCHMARKING REASONING LLMs

The development of a robust benchmark is crucial for documenting the advancements in reasoning LLMs capabilities and for identifying promising research directions for future progress. Here, we review the benchmarks from three key aspects: categories, evaluation metrics, and performance comparisons, while offering our reflections and insights.

4.1 Benchmark Categories

We categorize reasoning benchmarks by task type, which can be broadly divided into math, code, scientific, agent, medical, and multimodal reasoning. The detailed statistics for these benchmarks are presented in Table 6.

4.1.1 Benchmark Introduction

- 1) **Math Problems:** We document the current popular competition-level mathematical benchmarks to showcase the capabilities of reasoning LLMs, including AIME 2024 [226], MATH-500 [37], AMC 2023 [227], and Olympiad Bench [228].

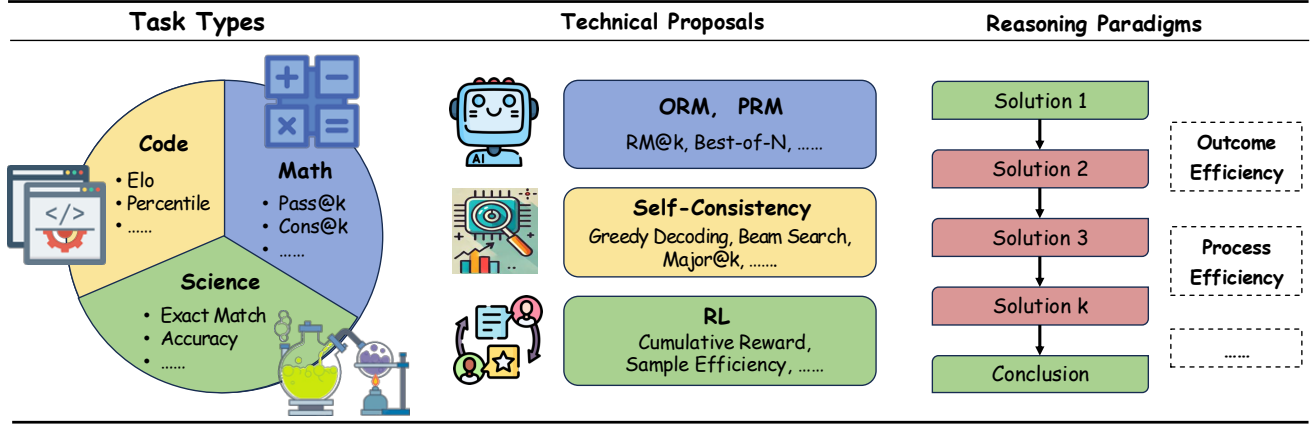


Fig. 6. Various evaluation metrics of reasoning LLMs divided by task types, technical proposals, and reasoning paradigms.

- 2) **Code Problems:** Code problems requires solid foundation and high logical thinking to evaluate the reasoning ability of reasoning LLMs such as Codeforces, SWE-bench [229], and LiveCodeBench [230].
- 3) **Scientific Problems:** Scientific benchmarks, *i.e.*, GPQA Diamond [231] and MMLU-Pro [232], involve multi-domains reasoning about chemistry, biology, and physics, which requires extensive knowledge accumulation and integrated reasoning.
- 4) **Agent Reasoning:** Realistic tasks often involve complex planning and tool usage, leading to the creation of agent reasoning benchmarks [244]. For example, WebShop [233] and WebArena [234] focus on web operations, while SciWorld [235] and TextCraft [236] are centered around scientific research.
- 5) **Medical Reasoning:** Medicine fundamentally involves complex reasoning, spanning tasks from diagnostic decision making to treatment planning. Benchmarks of JAMA Clinical Challenge [237], Medbullets [237], and MedQA [238] offer model measurements that mimic the doctor’s disease diagnosis.
- 6) **Multimodal Reasoning:** Multimodal reasoning, such as benchmarks of MMMU [239] and MathVista [240], requires cross-modal thinking in combination with text and images. Especially for those visual-centered problems, in benchmarks MathVision [241], MathVerse [245], CMMaTH [242], and PGPS9K [243], put forward higher requirements for reasoning LLMs.

4.1.2 Summary

The field of LLMs has advanced rapidly in recent years, with benchmark performance consistently improving. Simple reasoning benchmarks, such as GSM8K [32], MATH-500 [37], and ScienceQA [246], have approached performance saturation. Recent studies on reasoning LLMs [54], [147] show that models designed for long reasoning chains do not significantly outperform those designed for shorter chains on these benchmarks. This highlights the urgent need to establish new benchmarks that more effectively assess the reasoning capabilities of reasoning LLMs. Moreover, current benchmarks are limited, focusing mainly on solid reasoning tasks. Soft reasoning benchmarks, lacking explicitly defined correct answers, offer a more nuanced evaluation, better

capturing the complexities and subtleties of human-like reasoning. Furthermore, it is essential to address the issue of data leakage in evaluation processes [247]. Ensuring the confidentiality and neutrality of evaluation data is critical to preserving the integrity and reliability of benchmark results.

4.2 Evaluation Metrics

Depending on task types, technical proposals, and reasoning paradigms, various evaluation metrics have been introduced for reasoning LLMs as shown in Figure 6. These metrics are designed to more accurately assess the model’s performance in handling complex reasoning tasks, ensuring that both the quality and coherence of the generated solutions are effectively measured.

4.2.1 Task Types

In terms of benchmark categories, mathematical reasoning typically uses two main metrics: *Pass@k* and *Cons@k*. The *Pass@k* metric evaluates the model’s ability to generate a correct solution within *k* attempts, measuring the likelihood of success within a limited number of tries. On the other hand, *Cons@k* assesses whether the model consistently produces correct or logically coherent solutions, highlighting the stability and reliability of its reasoning capabilities. For code tasks, the key metrics are *Elo* and *Percentile*, both of which measure the relative skill in generating correct code compared to other models or human programmers. In scientific tasks, evaluation generally employs *Exact Match* (*EM*) and *Accuracy* for fill-in-the-blank and multiple-choice questions, respectively. The *EM* metric judges whether the model’s output exactly matches the expected solution, while *Accuracy* measures the proportion of correct answers out of the total number of questions.

4.2.2 Technical Proposals

Based on technical routes, the schemes with ORM or PRM often leverage *RM@k* and *Best-of-N* two evaluation indicators. *RM@k* measures whether the reward model can rank the good answer higher in the top *k* candidates according to reward score, and *Best-of-N* chooses the solution with highest score from *N* generated reasoning trajectories. Methods for self-consistency are evaluated using *Greedy Decoding*,

TABLE 7

Performance of Different Models, including Basic LLMs and Reasoning LLMs, on Plain Text Benchmarks. The **red** denotes the highest result, and the **blue** denotes the second highest result.

	Model	Math		Code			General	
		AIME 2024 (Pass@1)	MATH-500 (Pass@1)	LiveCodeBench (Pass@1-CoT)	Codeforces (Percentile)	SWE Verified (Resolved)	MMLU (Pass@1)	GPQA-Diamond (Pass@1)
Basic LLMs	GPT-4o [16]	9.3	74.6	34.2	23.6	38.8	87.2	49.9
	Claude-3.5-Sonnet [248]	16.0	78.3	33.8	20.3	50.8	88.3	65.0
	Gemini-2.0-Pro [249]	-	91.8	36.0	-	-	86.5	64.7
	Deepseek-V3 [17]	39.2	90.2	36.2	58.7	42.0	88.5	59.1
Reasoning LLMs	Eurus-2-7B-PRIME [195]	26.7	79.2	-	-	-	-	-
	InternLM3-8B-Instruct [250]	20.0	83.0	-	-	-	76.6	37.4
	rStar-Math-7B [147]	46.7	81.6	-	-	-	82.7	54.9
	STILL-2-32B [53]	46.7	90.2	-	-	-	-	-
	Redstar-code-math [54]	53.3	91.2	-	-	-	-	-
	Search-o1 [97]	56.7	86.4	33.0	-	-	-	63.6
	QwQ [98]	50.0	90.6	41.9	62.0	-	-	54.5
	s1-32B [251]	56.7	93.0	-	-	-	-	59.6
	OpenAI o1-mini [252]	63.6	90.0	53.8	93.4	41.6	85.2	60.0
	LIMO-32B [253]	57.1	94.8	-	-	-	-	66.7
	Kimi k1.5 long-CoT [192]	77.5	96.2	62.5	94.0	-	-	-
	DeepSeek-R1 [31]	79.8	97.3	65.9	96.3	49.2	90.8	71.5
	OpenAI-o1 [29]	79.2	96.4	63.4	96.6	48.9	91.8	75.7
	OpenAI o3-mini [30]	87.3	97.9	84.6	-	49.3	86.9	79.7

TABLE 8

Performance of Models, including Basic LLMs and Reasoning LLMs, on Multimodal Benchmarks. The **red** denotes the highest result, and the **blue** denotes the second highest result.

	Model	MMMU	Mathvista	Mathvision	Olympiadbench
Basic LLMs	GPT-4o [16]	69.1	63.8	30.4	25.9
	Claude-3.5-Sonnet [248]	70.4	65.3	35.6	-
	Gemini 2.0 Pro [249]	72.7	-	-	-
	LLaVA-CoT [176]	-	54.8	-	-
Reasoning LLMs	QvQ-72B-preview [254]	70.3	71.4	35.9	20.4
	Kimi k1.5 long-CoT [192]	70.0	74.9	-	-
	OpenAI-o1 [29]	77.3	71.0	-	-

Beam Search, and *Major@k*. *Greedy Decoding* and *Beam Search* control the randomness of the inference process by limiting the sampling range. *Major@k* selects the solution with the most consistent results from k candidate solutions. In RL, metrics reflect both performance in achieving desired outcomes and the efficiency of the learning process. For example, *Cumulative Reward* measures the total reward received by the agent over time, while *Sample Efficiency* assesses the efficiency of the agent’s sample usage during learning.

4.2.3 Reasoning Paradigms

For reasoning paradigm of the multi-turn solution generation in reasoning LLMs, *Outcome Efficiency* and *Process Efficiency* [102] are proposed recently to evaluate the efficiency of long thinking specifically. *Outcome Efficiency* metric empirically evaluates how effectively later solutions contribute to accuracy improvements, formulating as the ratio of efficient tokens that contribute to reaching the correct answer, to all output tokens. *Process Efficiency* metric evaluates the contribution of later solutions to solution diversity empirically, concretely representing as the ratio of tokens of distinct solutions to all solution tokens. These two indicators reveal to the overthinking issue of existing reasoning LLMs to simple problems certainly.

4.2.4 Summary

Most of the existing evaluation metrics are judged according to the final answer. It is imperative to develop a com-

prehensive assessment framework that considers various aspects of the reasoning process in view of the large inference computation consumption. Current popular evaluation frameworks, such as LMMs-Eval [255], OpenCompass [256], and PRMBench [257], lack efficiency and their metrics do not adequately account for the computational and temporal efficiency of the reasoning process. To address these shortcomings, we highly recommend exploring more efficient proxy tasks as potential solutions. By identifying and utilizing tasks that better capture the nuances of long reasoning chains, we can develop more robust and effective evaluation metrics to enhance the overall assessment framework, ensuring that it not only measures the accuracy of the final output but also evaluates the efficiency and coherence of the reasoning process throughout.

4.3 Performance Comparison

In this section, we compare the performance of different reasoning LLMs and their corresponding foundational LLMs on plain text benchmarks, such as math and code problems, as well as on multimodal benchmarks. The comprehensive real-time leaderboard is available on this website.

4.3.1 Performance on Plain Text Benchmarks

As shown in Table 7, reasoning LLMs, such as DeepSeek-R1 [31] and OpenAI-o1/o3 [29], [30], demonstrate exceptional performance across a wide range of tasks, including math, coding, and other general tasks. These models achieve high scores on multiple plain-text benchmarks, such as AIME 2024, MATH-500, and LiveCodeBench, showcasing their robust text-based reasoning abilities. In contrast, foundational LLMs, like GPT-4o [62], Claude-3.5-Sonnet [248], and DeepSeek-V3 [17], generally perform less effectively than reasoning LLMs, particularly in math and coding tasks (e.g., AIME 2024 and Codeforces). For example, OpenAI-o1 outperforms GPT-4o by 69.9% and 73% on these tasks, respectively. Moreover, DeepSeek-R1, based on the DeepSeek-V3 architecture, surpasses its predecessor on all benchmarks, further highlighting the advantages of the reasoning LLMs.

4.3.2 Performance on Multimodal Benchmarks

As shown in Table 8, reasoning LLMs continue to excel in multimodal tasks. OpenAI-o1 [29] performs strongly in vision tasks, achieving the highest score of 77.3% on MMMU and outperforming its corresponding foundational LLM, GPT-4o [62], by 7.2% on MathVista. However, the performance improvement in multimodal tasks is less pronounced compared to text-only tasks. This can be attributed in part to the limitations of current multimodal reasoning LLM techniques, as well as the lack of sufficient datasets to fully assess the multimodal capabilities of reasoning LLMs.

4.3.3 Summary

In summary, reasoning LLMs show strong performance across both plain text and multimodal benchmarks, particularly excelling in math and coding tasks, where they outperform foundational LLMs by a large margin. Although the improvement in multimodal tasks is not as pronounced as in text-only tasks, reasoning LLMs still surpass their counterparts, highlighting their potential for processing both image and text data. These results emphasize the versatility and effectiveness of reasoning LLMs across a broad spectrum of reasoning tasks, with potential for further advancements in multimodal reasoning techniques.

5 CHALLENGES & FUTURE DIRECTIONS

Despite the rapid advancements in reasoning LLMs, several challenges persist, limiting their generalizability and practical applicability. This section outlines these challenges and highlights potential research directions to address them.

5.1 Efficient Reasoning LLMs

While reasoning LLMs excel at solving complex problems via extended inference, their reliance on long autoregressive reasoning within large-scale architectures presents significant efficiency challenges. For example, many problems on platforms like Codeforces require over 10,000 tokens of reasoning, resulting in high latency. As noted in [102], even when a reasoning LLM identifies the correct solution early, it often spends considerable time verifying its reasoning. Recent reports, such as Deepseek-R1 [31], suggest that self-improvement via RL is more effective in larger models, while smaller-scale large language models (SLMs) (e.g., 3B and 7B models as explored by [103] and [199], [216]) struggle to match performance in slow-thinking reasoning tasks.

Future research should focus on two key areas: (1) integrating external reasoning tools to enable early stopping and verification mechanisms, thus improving the efficiency of long inference chains, and (2) exploring strategies to implement slow-thinking reasoning capabilities in SLMs without sacrificing performance.

5.2 Collaborative Slow & Fast-thinking Systems

A key challenge in reasoning LLMs is the loss of fast-thinking capabilities, which results in inefficiencies when simple tasks require unnecessary deep reasoning. Unlike humans, who fluidly switch between fast (*System 1*) and slow (*System 2*) thinking, current reasoning LLMs struggle to maintain this balance. While reasoning LLMs ensure

deliberate and thorough reasoning, fast-thinking systems rely on prior knowledge for quick responses. Despite efforts such as the *System 1-2* switcher [95], speculative decoding [258]–[260], and interactive continual learning [261], integrating both modes of thinking remains challenging. This often leads to inefficiencies in domain-specific tasks and underutilized strengths in more complex scenarios.

Future research should focus on developing adaptive switching mechanisms, joint training frameworks, and co-evolution strategies to harmonize the efficiency of fast-thinking systems with the precision of reasoning LLMs. Achieving this balance is crucial for advancing the field and creating more versatile AI systems.

5.3 Reasoning LLMs For Science

Reasoning LLMs play a crucial role in scientific research [262], enabling deep, structured analysis that goes beyond the heuristic-based fast-thinking models. Their value becomes especially clear in fields that demand complex reasoning, such as medicine and mathematics. In medicine, particularly in differential diagnosis and treatment planning, reasoning LLMs (e.g., inference-time scaling) enhance AI’s step-by-step reasoning, improving diagnostic accuracy where traditional scaling methods fall short [52]. In mathematics, approaches like FunSearch [263] incorporate slow-thinking principles to push beyond previous discoveries, showcasing the potential of AI-human collaboration.

Beyond these fields, reasoning LLMs can foster advancements in physics, engineering, and computational biology by refining model formulation and hypothesis testing. Investing in reasoning LLMs research not only bridges the gap between AI’s computational power and human-like analytical depth but also paves the way for more reliable, interpretable, and groundbreaking scientific discoveries.

5.4 Deep Integration of Neural and Symbolic Systems

Despite significant advancements in reasoning LLMs, their limited transparency and interpretability restrict their performance in more complex real-world reasoning tasks. The reliance on large-scale data patterns and lack of clear reasoning pathways makes it challenging to handle intricate or ambiguous problems effectively. Early symbolic logic systems, while less adaptable, offered better explainability and clearer reasoning steps, leading to more reliable performance in such cases.

A promising future direction is the deep integration of neural and symbolic systems. Google’s AlphaGeometry [264] and AlphaGeometry2 [265] combine reasoning LLMs with symbolic engines, achieving breakthroughs in the International Olympiad in Mathematics (IMO). In particular, AlphaGeometry2 utilizes the Gemini-based model [249], [266], [267] and a more efficient symbolic engine, improving performance by reducing rule sets and enhancing key concept handling. The system now covers a broader range of geometric concepts, including locus theorems and linear equations. A new search algorithm and knowledge-sharing mechanism accelerate the process. This system solved 84% of IMO geometry problems (2000-2024), surpassing gold medalists’ averages. In contrast, reasoning LLMs like OpenAI-o1 [29] failed to solve any problems. The

integration of neural and symbolic systems offers a balanced approach, improving both adaptability and interpretability, with vast potential for complex real-world reasoning tasks beyond mathematical geometry problems.

5.5 Multilingual Reasoning LLMs

Current reasoning LLMs perform well in high-resource languages like English and Chinese, demonstrating strong capabilities in tasks such as translation and various reasoning tasks [93], [101]. These models excel in environments where large-scale data and diverse linguistic resources are available. However, their performance in low-resource languages remains limited [268], facing challenges related to data sparsity, stability, safety, and overall performance. These issues hinder the effectiveness of reasoning LLMs in languages that lack substantial linguistic datasets and resources.

Future research should prioritize overcoming the challenges posed by data scarcity and cultural biases in low-resource languages. Innovations such as parameter sharing across reasoning LLMs and the incremental injection of domain-specific knowledge could help mitigate these challenges, enabling faster adaptation of slow-thinking capabilities to a broader range of languages. This would not only enhance the effectiveness of reasoning LLMs in these languages but also ensure more equitable access to advanced AI technologies.

5.6 Multimodal Reasoning LLMs

Extending slow-thinking reasoning capabilities from text-based domains to multimodal contexts remains a significant challenge, especially in tasks requiring fine-grained perception [96]. While approaches like Virgo [269] have attempted to distill text-based slow-thinking reasoning into multimodal LLMs, their performance improvements in tasks such as MathVision [241], which demand detailed visual understanding, have been marginal.

Key research directions include developing hierarchical reasoning LLMs that enable fine-grained cross-modal understanding and generation, tailored to the unique characteristics of modalities such as audio, video, and 3D data.

5.7 Safe Reasoning LLMs

The rapid development of reasoning LLMs like OpenAI-o1 [29] and DeepSeek-R1 [31] has led to the rise of superintelligent models capable of continuous self-evolution. However, this progress brings challenges in safety and control. RL, a key training method, introduces risks such as reward hacking, generalization failures, and language mixing, which can lead to harmful outcomes. Ensuring the safety of such systems like DeepSeek-R1 is urgent. While RL enhances reasoning, its uncontrollable nature raises concerns about safely guiding these models. SFT addresses some issues but is not a complete solution. A hybrid approach combining RL and SFT is needed to reduce harmful outputs while maintaining model effectiveness [270].

As these models surpass human cognitive capabilities, ensuring their safe, responsible, and transparent use is crucial. This requires ongoing research to develop methods for controlling and guiding their actions, thereby balancing AI power with ethical decision-making.

6 CONCLUSION

This paper presents a comprehensive survey that advances research on reasoning LLMs. We begin with an overview of the progress in foundational LLMs and key early *System 2* technologies, including symbolic logic, MCTS, and RL, exploring how each, when combined with foundational LLMs, has paved the way for reasoning LLMs. We then provide a detailed feature analysis of the latest reasoning LLMs, examining the core methods that enable their advanced reasoning capabilities and highlighting representative models. Through a review of mainstream reasoning benchmarks and performance comparisons, we offer valuable insights into the current state of the field. Looking ahead, we identify promising research directions and continue to track developments via our real-time GitHub Repository. This survey aims to inspire innovation and foster progress in the rapidly evolving field of reasoning LLMs.

REFERENCES

- [1] W. Hua and Y. Zhang, "System 1+ system 2= better world: Neural-symbolic chain of logic reasoning," in *Findings of the Association for Computational Linguistics: EMNLP 2022*, 2022, pp. 601–612.
- [2] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Advances in neural information processing systems*, vol. 35, pp. 24 824–24 837, 2022.
- [3] X. Wang, J. Wei, D. Schuurmans, Q. V. Le, E. H. Chi, S. Narang, A. Chowdhery, and D. Zhou, "Self-Consistency Improves Chain of Thought Reasoning in Language Models," in *The Eleventh International Conference on Learning Representations*, 2023.
- [4] D. Zhou, N. Schärli, L. Hou, J. Wei, N. Scales, X. Wang, D. Schuurmans, C. Cui, O. Bousquet, Q. V. Le *et al.*, "Least-to-Most Prompting Enables Complex Reasoning in Large Language Models," in *The Eleventh International Conference on Learning Representations*, 2023.
- [5] E. Zelikman, Y. Wu, J. Mu, and N. D. Goodman, "STaR: Self-taught reasoner bootstrapping reasoning with reasoning," in *Proc. the 36th International Conference on Neural Information Processing Systems*, vol. 1126, 2024.
- [6] J. S. B. Evans, "Heuristic and analytic processes in reasoning," *British Journal of Psychology*, vol. 75, no. 4, pp. 451–468, 1984.
- [7] D. Kahneman, "Maps of bounded rationality: Psychology for behavioral economics," *American economic review*, vol. 93, no. 5, pp. 1449–1475, 2003.
- [8] J. Huang and K. C.-C. Chang, "Towards Reasoning in Large Language Models: A Survey," in *Findings of the Association for Computational Linguistics: ACL 2023*, 2023, pp. 1049–1065.
- [9] S. Qiao, Y. Ou, N. Zhang, X. Chen, Y. Yao, S. Deng, C. Tan, F. Huang, and H. Chen, "Reasoning with Language Model Prompting: A Survey," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2023, pp. 5368–5393.
- [10] B. Wang, S. Min, X. Deng, J. Shen, Y. Wu, L. Zettlemoyer, and H. Sun, "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)*, 2023, pp. 2717–2739.
- [11] O. Shaikh, H. Zhang, W. Held, M. Bernstein, and D. Yang, "On Second Thought, Let's Not Think Step by Step! Bias and Toxicity in Zero-Shot Reasoning," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2023, pp. 4454–4470.
- [12] H. Shao, S. Qian, H. Xiao, G. Song, Z. Zong, L. Wang, Y. Liu, and H. Li, "Visual cot: Advancing multi-modal language models with a comprehensive dataset and benchmark for chain-of-thought reasoning," in *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
- [13] Z. Zhang, A. Zhang, M. Li, and A. Smola, "Automatic Chain of Thought Prompting in Large Language Models," in *The Eleventh International Conference on Learning Representations*, 2023.

- [14] S. Hao, Y. Gu, H. Ma, J. Hong, Z. Wang, D. Wang, and Z. Hu, "Reasoning with Language Model is Planning with World Model," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023, pp. 8154–8173.
- [15] Y. Zhang, "Meta prompting for agi systems," *arXiv preprint arXiv:2311.11482*, 2023.
- [16] OpenAI, "Hello GPT-4o," May 2024. [Online]. Available: <https://openai.com/index/hello-gpt-4o/>
- [17] A. Liu, B. Feng, B. Xue, B. Wang, B. Wu, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan *et al.*, "Deepseek-v3 technical report," *arXiv preprint arXiv:2412.19437*, 2024.
- [18] A. Vaswani, "Attention is all you need," *Advances in Neural Information Processing Systems*, 2017.
- [19] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, 2019, pp. 4171–4186.
- [20] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "RoBERTa: A Robustly Optimized BERT Pretraining Approach," *CoRR*, vol. abs/1907.11692, 2019.
- [21] A. Radford, "Improving language understanding by generative pre-training," 2018.
- [22] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [23] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [24] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, "Training language models to follow instructions with human feedback," *Advances in neural information processing systems*, vol. 35, pp. 27 730–27 744, 2022.
- [25] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.
- [26] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, "A survey of large language models," *arXiv preprint arXiv:2303.18223*, 2023.
- [27] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual Instruction Tuning," in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [28] D. Zhang, Y. Yu, J. Dong, C. Li, D. Su, C. Chu, and D. Yu, "MM-LLMs: Recent Advances in MultiModal Large Language Models," in *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*. Association for Computational Linguistics, 2024, pp. 12 401–12 430.
- [29] OpenAI, "Learning to reason with LLMs," September 2024. [Online]. Available: <https://openai.com/index/learning-to-reason-with-llms/>
- [30] —, "OpenAI o3-mini," January 2025. [Online]. Available: <https://openai.com/index/openai-o3-mini/>
- [31] D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi *et al.*, "DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning," *arXiv preprint arXiv:2501.12948*, 2025.
- [32] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano *et al.*, "Training verifiers to solve math word problems," *arXiv preprint arXiv:2110.14168*, 2021.
- [33] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, "Large language models are zero-shot reasoners," *Advances in neural information processing systems*, vol. 35, pp. 22 199–22 213, 2022.
- [34] Y. Liu, A. Singh, C. D. Freeman, J. D. Co-Reyes, and P. J. Liu, "Improving large language model fine-tuning for solving math problems," *arXiv preprint arXiv:2310.10047*, 2023.
- [35] X. Zhu, J. Wang, L. Zhang, Y. Zhang, Y. Huang, R. Gan, J. Zhang, and Y. Yang, "Solving Math Word Problems via Cooperative Reasoning induced Language Models," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2023, pp. 4471–4485.
- [36] P. Lu, L. Qiu, K.-W. Chang, Y. N. Wu, S.-C. Zhu, T. Rajpurohit, P. Clark, and A. Kalyan, "Dynamic Prompt Learning via Policy Gradient for Semi-structured Mathematical Reasoning," in *The Eleventh International Conference on Learning Representations*, 2023.
- [37] H. Lightman, V. Kosaraju, Y. Burda, H. Edwards, B. Baker, T. Lee, J. Leike, J. Schulman, I. Sutskever, and K. Cobbe, "Let's Verify Step by Step," in *The Twelfth International Conference on Learning Representations*, 2024.
- [38] F. Yao, C. Tian, J. Liu, Z. Zhang, Q. Liu, L. Jin, S. Li, X. Li, and X. Sun, "Thinking like an expert: Multimodal hypergraph-of-thought (hot) reasoning to boost foundation modals," *arXiv preprint arXiv:2308.06207*, 2023.
- [39] Y. Yao, Z. Li, and H. Zhao, "Beyond Chain-of-Thought, Effective Graph-of-Thought Reasoning in Language Models," *arXiv preprint arXiv:2305.16582*, 2023.
- [40] Y. Wen, Z. Wang, and J. Sun, "Mindmap: Knowledge graph prompting sparks graph of thoughts in large language models," *arXiv preprint arXiv:2308.09729*, 2023.
- [41] B. Lei, C. Liao, C. Ding *et al.*, "Boosting logical reasoning in large language models through a new framework: The graph of thought," *arXiv preprint arXiv:2308.08614*, 2023.
- [42] M. Jin, Q. Yu, D. Shu, H. Zhao, W. Hua, Y. Meng, Y. Zhang, and M. Du, "The impact of reasoning step length on large language models," *arXiv preprint arXiv:2401.04925*, 2024.
- [43] M. Besta, N. Blach, A. Kubicek, R. Gerstenberger, M. Podstawski, L. Gianinazzi, J. Gajda, T. Lehmann, H. Niewiadomski, P. Nyczzyk *et al.*, "Graph of thoughts: Solving elaborate problems with large language models," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 16, 2024, pp. 17 682–17 690.
- [44] P. Cheng, T. Hu, H. Xu, Z. Zhang, Y. Dai, L. Han, and N. Du, "Self-playing Adversarial Language Game Enhances LLM Reasoning," *arXiv preprint arXiv:2404.10642*, 2024.
- [45] H. You, R. Sun, Z. Wang, L. Chen, G. Wang, H. Ayyubi, K.-W. Chang, and S.-F. Chang, "IdealGPT: Iteratively Decomposing Vision and Language Reasoning via Large Language Models," in *Findings of the Association for Computational Linguistics: EMNLP 2023*, 2023, pp. 11 289–11 303.
- [46] P. Wu and S. Xie, "V?: Guided Visual Search as a Core Mechanism in Multimodal LLMs," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 13 084–13 094.
- [47] Z. Chen, R. Sun, W. Liu, Y. Hong, and C. Gan, "GENOME: Generative Neuro-Symbolic Visual Reasoning by Growing and Reusing Modules," in *International Conference on Learning Representations*, 2024.
- [48] S. Wu, Z. Peng, X. Du, T. Zheng, M. Liu, J. Wu, J. Ma, Y. Li, J. Yang, W. Zhou *et al.*, "A Comparative Study on Reasoning Patterns of OpenAI's o1 Model," *arXiv preprint arXiv:2410.13639*, 2024.
- [49] V. Xiang, C. Snell, K. Gandhi, A. Albalak, A. Singh, C. Blagden, D. Phung, R. Rafailov, N. Lile, D. Mahan *et al.*, "Towards System 2 Reasoning in LLMs: Learning How to Think With Meta Chain-of-Thought," *arXiv preprint arXiv:2501.04682*, 2025.
- [50] Y. Qin, X. Li, H. Zou, Y. Liu, S. Xia, Z. Huang, Y. Ye, W. Yuan, H. Liu, Y. Li *et al.*, "O1 Replication Journey: A Strategic Progress Report—Part 1," *arXiv preprint arXiv:2410.18982*, 2024.
- [51] Z. Huang, H. Zou, X. Li, Y. Liu, Y. Zheng, E. Chern, S. Xia, Y. Qin, W. Yuan, and P. Liu, "O1 Replication Journey—Part 2: Surpassing O1-preview through Simple Distillation, Big Progress or Bitter Lesson?" *arXiv preprint arXiv:2411.16489*, 2024.
- [52] Z. Huang, G. Geng, S. Hua, Z. Huang, H. Zou, S. Zhang, P. Liu, and X. Zhang, "O1 Replication Journey—Part 3: Inference-time Scaling for Medical Reasoning," *arXiv preprint arXiv:2501.06458*, 2025.
- [53] Y. Min, Z. Chen, J. Jiang, J. Chen, J. Deng, Y. Hu, Y. Tang, J. Wang, X. Cheng, H. Song, W. X. Zhao, Z. Liu, Z. Wang, and J.-R. Wen, "Imitate, Explore, and Self-Improve: A Reproduction Report on Slow-thinking Reasoning Systems," *arXiv preprint arXiv:2412.09413*, 2024.
- [54] H. Xu, X. Wu, W. Wang, Z. Li, D. Zheng, B. Chen, Y. Hu, S. Kang, J. Ji, Y. Zhang *et al.*, "RedStar: Does Scaling Long-CoT Data Unlock Better Slow-Reasoning Systems?" *arXiv preprint arXiv:2501.11284*, 2025.
- [55] Z. Zeng, Q. Cheng, Z. Yin, B. Wang, S. Li, Y. Zhou, Q. Guo, X. Huang, and X. Qiu, "Scaling of Search and Learning: A

- Roadmap to Reproduce o1 from Reinforcement Learning Perspective," *arXiv preprint arXiv:2412.14135*, 2024.
- [56] Y. Ji, J. Li, H. Ye, K. Wu, J. Xu, L. Mo, and M. Zhang, "Test-time Computing: from System-1 Thinking to System-2 Thinking," *arXiv preprint arXiv:2501.02497*, 2025.
- [57] M. Besta, J. Barth, E. Schreiber, A. Kubicek, A. Catarino, R. Gerstenberger, P. Nyczyk, P. Iff, Y. Li, S. Houliston *et al.*, "Reasoning Language Models: A Blueprint," *arXiv preprint arXiv:2501.11223*, 2025.
- [58] Y. Zhang, S. Mao, T. Ge, X. Wang, A. de Wuyter, Y. Xia, W. Wu, T. Song, M. Lan, and F. Wei, "LLM as a Mastermind: A Survey of Strategic Reasoning with Large Language Models," *arXiv preprint arXiv:2404.01230*, 2024.
- [59] F. Xu, Q. Hao, Z. Zong, J. Wang, Y. Zhang, J. Wang, X. Lan, J. Gong, T. Ouyang, F. Meng *et al.*, "Towards Large Reasoning Models: A Survey of Reinforced Reasoning with Large Language Models," *arXiv preprint arXiv:2501.09686*, 2025.
- [60] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [61] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, "Zero-shot text-to-image generation," in *International conference on machine learning*. Pmlr, 2021, pp. 8821–8831.
- [62] OpenAI, "GPT-4 Technical Report," 2023.
- [63] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural computation*, vol. 18, no. 7, pp. 1527–1554, 2006.
- [64] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath *et al.*, "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *IEEE Signal processing magazine*, vol. 29, no. 6, pp. 82–97, 2012.
- [65] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, 2012.
- [66] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [67] Q. Dong, L. Li, D. Dai, C. Zheng, J. Ma, R. Li, H. Xia, J. Xu, Z. Wu, B. Chang *et al.*, "A survey on in-context learning," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024, pp. 1107–1128.
- [68] J. White, Q. Fu, S. Hays, M. Sandborn, C. Olea, H. Gilbert, A. Elnashar, J. Spencer-Smith, and D. C. Schmidt, "A prompt pattern catalog to enhance prompt engineering with chatgpt," *arXiv preprint arXiv:2302.11382*, 2023.
- [69] C. I. Lewis, C. H. Langford, and P. Lamprecht, *Symbolic logic*. Dover publications New York, 1959, vol. 170.
- [70] R. Carnap, *Introduction to symbolic logic and its applications*. Courier Corporation, 2012.
- [71] A. Colmerauer, "An introduction to Prolog III," *Communications of the ACM*, vol. 33, no. 7, pp. 69–90, 1990.
- [72] W. F. Clocksin and C. S. Mellish, *Programming in PROLOG*. Springer Science & Business Media, 2003.
- [73] K. R. Apt *et al.*, *From logic programming to Prolog*. Prentice Hall London, 1997, vol. 362.
- [74] M. P. Singh, A. S. Rao, and M. P. Georgeff, *Formal methods in DAI: Logic-based representation and reasoning*. MIT Press Cambridge, 1999, vol. 8.
- [75] R. G. Jeroslow, "Computation-oriented reductions of predicate to propositional logic," *Decision Support Systems*, vol. 4, no. 2, pp. 183–197, 1988.
- [76] J. McCarthy, "History of LISP," in *History of programming languages*, 1978, pp. 173–185.
- [77] L. Bachmair and H. Ganzinger, "Resolution Theorem Proving," *Handbook of automated reasoning*, vol. 1, no. 02, 2001.
- [78] M. Minsky *et al.*, "A framework for representing knowledge," 1974.
- [79] C. B. Browne, E. Powley, D. Whitehouse, S. M. Lucas, P. I. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, and S. Colton, "A survey of monte carlo tree search methods," *IEEE Transactions on Computational Intelligence and AI in games*, vol. 4, no. 1, pp. 1–43, 2012.
- [80] S. Gelly and D. Silver, "Monte-Carlo tree search and rapid action value estimation in computer Go," *Artificial Intelligence*, vol. 175, no. 11, pp. 1856–1875, 2011.
- [81] M. Świechowski, K. Godlewski, B. Sawicki, and J. Mańdziuk, "Monte Carlo tree search: A review of recent modifications and applications," *Artificial Intelligence Review*, vol. 56, no. 3, pp. 2497–2562, 2023.
- [82] R. S. Sutton, A. G. Barto *et al.*, *Reinforcement learning: An introduction*. MIT press Cambridge, 1998, vol. 1, no. 1.
- [83] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, pp. 279–292, 1992.
- [84] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [85] R. R. Torrado, P. Bontrager, J. Togelius, J. Liu, and D. Perez-Liebana, "Deep reinforcement learning for general video game ai," in *2018 IEEE Conference on Computational Intelligence and Games (CIG)*. IEEE, 2018, pp. 1–8.
- [86] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, "Mastering the game of Go with deep neural networks and tree search," *nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [87] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton *et al.*, "Mastering the game of go without human knowledge," *nature*, vol. 550, no. 7676, pp. 354–359, 2017.
- [88] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev *et al.*, "Grandmaster level in StarCraft II using multi-agent reinforcement learning," *nature*, vol. 575, no. 7782, pp. 350–354, 2019.
- [89] A. Y. Ng, D. Harada, and S. Russell, "Policy invariance under reward transformations: Theory and application to reward shaping," in *icml*, vol. 99. Citeseer, 1999, pp. 278–287.
- [90] H. Luo, Q. Sun, C. Xu, P. Zhao, J. Lou, C. Tao, X. Geng, Q. Lin, S. Chen, and D. Zhang, "Wizardmath: Empowering mathematical reasoning for large language models via reinforced evolution-instruct," *arXiv preprint arXiv:2308.09583*, 2023.
- [91] Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, X. Bi, H. Zhang, M. Zhang, Y. Li, Y. Wu *et al.*, "Deepseekmath: Pushing the limits of mathematical reasoning in open language models," *arXiv preprint arXiv:2402.03300*, 2024.
- [92] E. Zelikman, G. Harik, Y. Shao, V. Jayasiri, N. Haber, and N. D. Goodman, "Quiet-star: Language models can teach themselves to think before speaking," *arXiv preprint arXiv:2403.09629*, 2024.
- [93] Y. Zhao, H. Yin, B. Zeng, H. Wang, T. Shi, C. Lyu, L. Wang, W. Luo, and K. Zhang, "Marco-o1: Towards open reasoning models for open-ended solutions," *arXiv preprint arXiv:2411.14405*, 2024.
- [94] J. Chen, Z. Cai, K. Ji, X. Wang, W. Liu, R. Wang, J. Hou, and B. Wang, "Huatuoqpt-o1, towards medical complex reasoning with llms," *arXiv preprint arXiv:2412.18925*, 2024.
- [95] G. Sun, M. Jin, Z. Wang, C.-L. Wang, S. Ma, Q. Wang, Y. N. Wu, Y. Zhang, and D. Liu, "Visual agents as fast and slow thinkers," *arXiv preprint arXiv:2408.08862*, 2024.
- [96] H. Wei, Y. Yin, Y. Li, J. Wang, L. Zhao, J. Sun, Z. Ge, and X. Zhang, "Slow Perception: Let's Perceive Geometric Figures Step-by-step," *arXiv preprint arXiv:2412.20631*, 2024.
- [97] X. Li, G. Dong, J. Jin, Y. Zhang, Y. Zhou, Y. Zhu, P. Zhang, and Z. Dou, "Search-o1: Agentic search-enhanced large reasoning models," *arXiv preprint arXiv:2501.05366*, 2025.
- [98] Q. Team, "QwQ: Reflect Deeply on the Boundaries of the Unknown," November 2024. [Online]. Available: <https://qwenlm.github.io/blog/qwq-32b-preview/>
- [99] N. Team, "Sky-T1: Train your own O1 preview model within \$450," <https://novasky-ai.github.io/posts/sky-t1>, 2025, accessed: 2025-01-09.
- [100] Y. Zhang, S. Wu, Y. Yang, J. Shu, J. Xiao, C. Kong, and J. Sang, "o1-coder: an o1 replication for coding," *arXiv preprint arXiv:2412.00154*, 2024.
- [101] J. Wang, F. Meng, Y. Liang, and J. Zhou, "DRT-o1: Optimized Deep Reasoning Translation via Long Chain-of-Thought," *arXiv preprint arXiv:2412.17498*, 2024.
- [102] X. Chen, J. Xu, T. Liang, Z. He, J. Pang, D. Yu, L. Song, Q. Liu, M. Zhou, Z. Zhang *et al.*, "Do NOT Think That Much for 2

- + 3=? On the Overthinking of o1-Like LLMs," *arXiv preprint arXiv:2412.21187*, 2024.
- [103] W. Zeng, Y. Huang, W. Liu, K. He, Q. Liu, Z. Ma, and J. He, "7B Model and 8K Examples: Emerging Reasoning with Reinforcement Learning is Both Effective and Efficient," <https://hkust-nlp.notion.site/simplerl-reason>, 2025, notion Blog.
- [104] Z. Wan, X. Feng, M. Wen, S. M. McAleer, Y. Wen, W. Zhang, and J. Wang, "Alphazero-like tree-search can guide large language model decoding and training," in *Forty-first International Conference on Machine Learning*, 2024.
- [105] Z. Bi, K. Han, C. Liu, Y. Tang, and Y. Wang, "Forest-of-Thought: Scaling Test-Time Compute for Enhancing LLM Reasoning," *CoRR*, vol. abs/2412.09078, 2024.
- [106] J. Li, H. Le, Y. Zhou, C. Xiong, S. Savarese, and D. Sahoo, "CodeTree: Agent-guided Tree Search for Code Generation with Large Language Models," *CoRR*, vol. abs/2411.04329, 2024.
- [107] J. Cheng, X. Liu, C. Wang, X. Gu, Y. Lu, D. Zhang, Y. Dong, J. Tang, H. Wang, and M. Huang, "SPaR: Self-Play with Tree-Search Refinement to Improve Instruction-Following in Large Language Models," *CoRR*, vol. abs/2412.11605, 2024.
- [108] J. Qiu, Y. Lu, Y. Zeng, J. Guo, J. Geng, H. Wang, K. Huang, Y. Wu, and M. Wang, "TreeBoN: Enhancing Inference-Time Alignment with Speculative Tree-Search and Best-of-N Sampling," *CoRR*, vol. abs/2410.16033, 2024.
- [109] N. Dainese, M. Merler, M. Alakuijala, and P. Marttinen, "Generating Code World Models with Large Language Models Guided by Monte Carlo Tree Search," *CoRR*, vol. abs/2405.15383, 2024.
- [110] Z. Zhao, W. S. Lee, and D. Hsu, "Large Language Models as Commonsense Knowledge for Large-Scale Task Planning," in *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [111] Q. Li, W. Xia, K. Du, X. Dai, R. Tang, Y. Wang, Y. Yu, and W. Zhang, "RethinkMCTS: Refining Erroneous Thoughts in Monte Carlo Tree Search for Code Generation," *CoRR*, vol. abs/2409.09584, 2024.
- [112] D. Zhang, X. Huang, D. Zhou, Y. Li, and W. Ouyang, "Accessing GPT-4 level Mathematical Olympiad Solutions via Monte Carlo Tree Self-refine with LLaMa-3 8B," *CoRR*, vol. abs/2406.07394, 2024.
- [113] G. Rabby, F. Keya, P. Zamil, and S. Auer, "MC-NEST – Enhancing Mathematical Reasoning in Large Language Models with a Monte Carlo Nash Equilibrium Self-Refine Tree," 2024. [Online]. Available: <https://arxiv.org/abs/2411.15645>
- [114] B. Xu, Y. Lin, Y. Li, and Y. Gao, "SRA-MCTS: Self-driven Reasoning Augmentation with Monte Carlo Tree Search for Code Generation," *CoRR*, vol. abs/2411.11053, 2024.
- [115] J. Kang, X. Z. Li, X. Chen, A. Kazemi, and B. Chen, "MindStar: Enhancing Math Reasoning in Pre-trained LLMs at Inference Time," *CoRR*, vol. abs/2405.16265, 2024.
- [116] P. Kadam, "GPT-Guided Monte Carlo Tree Search for Symbolic Regression in Financial Fraud Detection," *CoRR*, vol. abs/2411.04459, 2024.
- [117] D. Zhang, J. Wu, J. Lei, T. Che, J. Li, T. Xie, X. Huang, S. Zhang, M. Pavone, Y. Li, W. Ouyang, and D. Zhou, "LLaMA-Berry: Pairwise Optimization for O1-like Olympiad-Level Mathematical Reasoning," *CoRR*, vol. abs/2410.02884, 2024.
- [118] D. Zhang, S. Zhoubian, Y. Yue, Y. Dong, and J. Tang, "ReST-MCTS*: LLM Self-Training via Process Reward Guided Tree Search," *CoRR*, vol. abs/2406.03816, 2024.
- [119] Y. Xie, A. Goyal, W. Zheng, M. Kan, T. P. Lillicrap, K. Kawaguchi, and M. Shieh, "Monte Carlo Tree Search Boosts Reasoning via Iterative Preference Learning," *CoRR*, vol. abs/2405.00451, 2024.
- [120] G. Rabby, F. Keya, P. Zamil, and S. Auer, "MC-NEST - Enhancing Mathematical Reasoning in Large Language Models with a Monte Carlo Nash Equilibrium Self-Refine Tree," *CoRR*, vol. abs/2411.15645, 2024.
- [121] J. Y. Koh, S. McAleer, D. Fried, and R. Salakhutdinov, "Tree Search for Language Model Agents," *CoRR*, vol. abs/2407.01476, 2024.
- [122] J. Liu, A. Cohen, R. Pasunuru, Y. Choi, H. Hajishirzi, and A. Celikyilmaz, "Don't throw away your value model! Generating more preferable text with Value-Guided Monte-Carlo Tree Search decoding," 2024. [Online]. Available: <https://arxiv.org/abs/2309.15028>
- [123] C. Zhang, J. Song, S. Li, Y. Liang, Y. Ma, W. Wang, Y. Zhu, and S. Zhu, "Proposing and solving olympiad geometry with guided tree search," *CoRR*, vol. abs/2412.10673, 2024.
- [124] H. Jiang, Y. Ma, C. Ding, K. Luan, and X. Di, "Towards Intrinsic Self-Correction Enhancement in Monte Carlo Tree Search Boosted Reasoning via Iterative Preference Learning," 2024. [Online]. Available: <https://arxiv.org/abs/2412.17397>
- [125] H. Xu, "No Train Still Gain. Unleash Mathematical Reasoning of Large Language Models with Monte Carlo Tree Search Guided by Energy Function," *CoRR*, vol. abs/2309.03224, 2023.
- [126] M. Kemmerling, D. Lütticke, and R. H. Schmitt, "Beyond games: a systematic review of neural Monte Carlo tree search applications," *Appl. Intell.*, vol. 54, no. 11-12, pp. 1020–1046, 2024.
- [127] Y. Li, Z. Lin, S. Zhang, Q. Fu, B. Chen, J.-G. Lou, and W. Chen, "Making large language models better reasoners with step-aware verifier," *arXiv preprint arXiv:2206.02336*, 2022.
- [128] P. Wang, L. Li, Z. Shao, R. Xu, D. Dai, Y. Li, D. Chen, Y. Wu, and Z. Sui, "Math-shepherd: Verify and reinforce llms step-by-step without human annotations," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2024, pp. 9426–9439.
- [129] J. Lu, Z. Dou, W. Hongru, Z. Cao, J. Dai, Y. Feng, and Z. Guo, "Autopsv: Automated process-supervised verifier," in *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [130] L. Yuan, W. Li, H. Chen, G. Cui, N. Ding, K. Zhang, B. Zhou, Z. Liu, and H. Peng, "Free process rewards without process labels," *arXiv preprint arXiv:2412.01981*, 2024.
- [131] F. Yu, A. Gao, and B. Wang, "OVM, Outcome-supervised Value Models for Planning in Mathematical Reasoning," in *Findings of the Association for Computational Linguistics: NAACL 2024*, 2024, pp. 858–875.
- [132] D. Zhang, S. Zhoubian, Z. Hu, Y. Yue, Y. Dong, and J. Tang, "Rest-mcts*: Llm self-training via process reward guided tree search," *arXiv preprint arXiv:2406.03816*, 2024.
- [133] L. Luo, Y. Liu, R. Liu, S. Phatale, H. Lara, Y. Li, L. Shu, Y. Zhu, L. Meng, J. Sun *et al.*, "Improve Mathematical Reasoning in Language Models by Automated Process Supervision," *arXiv preprint arXiv:2406.06592*, 2024.
- [134] Z. Sun, Q. Wang, W. Yu, X. Zang, K. Zheng, J. Xu, X. Zhang, S. Yang, and H. Li, "ReARTeR: Retrieval-Augmented Reasoning with Trustworthy Process Rewarding," *arXiv preprint arXiv:2501.07861*, 2025.
- [135] Z. Zhang, C. Zheng, Y. Wu, B. Zhang, R. Lin, B. Yu, D. Liu, J. Zhou, and J. Lin, "The lessons of developing process reward models in mathematical reasoning," *arXiv preprint arXiv:2501.07301*, 2025.
- [136] Z. Yu, W. Gu, Y. Wang, Z. Zeng, J. Wang, W. Ye, and S. Zhang, "Outcome-Refining Process Supervision for Code Generation," *arXiv preprint arXiv:2412.15118*, 2024.
- [137] X. Lai, Z. Tian, Y. Chen, S. Yang, X. Peng, and J. Jia, "Step-dpo: Step-wise preference optimization for long-chain reasoning of llms," *arXiv preprint arXiv:2406.18629*, 2024.
- [138] Y. Liu, J. Lu, Z. Chen, C. Qu, J. K. Liu, C. Liu, Z. Cai, Y. Xia, L. Zhao, J. Bian *et al.*, "AdaptiveStep: Automatically Dividing Reasoning Step through Model Confidence," *arXiv preprint arXiv:2502.13943*, 2025.
- [139] F. Yu, A. Gao, and B. Wang, "Outcome-supervised verifiers for planning in mathematical reasoning," *arXiv preprint arXiv:2311.09724*, 2023.
- [140] J. Uesato, N. Kushman, R. Kumar, F. Song, N. Siegel, L. Wang, A. Creswell, G. Irving, and I. Higgins, "Solving math word problems with process-and outcome-based feedback," *arXiv preprint arXiv:2211.14275*, 2022.
- [141] Y. Li, Z. Lin, S. Zhang, Q. Fu, B. Chen, J.-G. Lou, and W. Chen, "Making language models better reasoners with step-aware verifier," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2023, pp. 5315–5333.
- [142] Z. Wu, Y. Hu, W. Shi, N. Dziri, A. Suhr, P. Ammanabrolu, N. A. Smith, M. Ostendorf, and H. Hajishirzi, "Fine-grained human feedback gives better rewards for language model training," *Advances in Neural Information Processing Systems*, vol. 36, pp. 59 008–59 033, 2023.
- [143] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn, "Direct preference optimization: Your language model is secretly a reward model," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [144] E. Zelikman, Y. Wu, J. Mu, and N. D. Goodman, "STaR: Bootstrapping Reasoning With Reasoning," in *Advances in Neural*

- Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022, 2022.*
- [145] A. Hosseini, X. Yuan, N. Malkin, A. Courville, A. Sordoni, and R. Agarwal, “V-STaR: Training Verifiers for Self-Taught Reasoners,” 2024. [Online]. Available: <https://arxiv.org/abs/2402.06457>
 - [146] W. Zeng, Y. Huang, L. Zhao, Y. Wang, Z. Shan, and J. He, “B-STaR: Monitoring and Balancing Exploration and Exploitation in Self-Taught Reasoners,” 2024. [Online]. Available: <https://arxiv.org/abs/2412.17256>
 - [147] X. Guan, L. L. Zhang, Y. Liu, N. Shang, Y. Sun, Y. Zhu, F. Yang, and M. Yang, “rStar-Math: Small LLMs Can Master Math Reasoning with Self-Evolved Deep Thinking,” 2025. [Online]. Available: <https://arxiv.org/abs/2501.04519>
 - [148] Ç. Gülçehre, T. L. Paine, S. Srinivasan, K. Konyushkova, L. Weerts, A. Sharma, A. Siddhant, A. Ahern, M. Wang, C. Gu, W. Macherey, A. Doucet, O. Firat, and N. de Freitas, “Reinforced Self-Training (ReST) for Language Modeling,” *CoRR*, vol. abs/2308.08998, 2023.
 - [149] A. Singh, J. D. Co-Reyes, R. Agarwal, A. Anand, P. Patil, X. Garcia, P. J. Liu, J. Harrison, J. Lee, K. Xu, A. Parisi, A. Kumar, A. Alemi, A. Rizkowsky, A. Nova, B. Adlam, B. Bohnet, G. Elsayed, H. Sedghi, I. Mordatch, I. Simpson, I. Gur, J. Snoek, J. Pennington, J. Hron, K. Kenealy, K. Swersky, K. Mahajan, L. Culp, L. Xiao, M. L. Bileschi, N. Constant, R. Novak, R. Liu, T. Warkentin, Y. Qian, Y. Bansal, E. Dyer, B. Neyshabur, J. Sohl-Dickstein, and N. Fiedel, “Beyond Human Data: Scaling Self-Training for Problem-Solving with Language Models,” 2024. [Online]. Available: <https://arxiv.org/abs/2312.06585>
 - [150] F. Xu, Q. Sun, K. Cheng, J. Liu, Y. Qiao, and Z. Wu, “Interactive Evolution: A Neural-Symbolic Self-Training Framework For Large Language Models,” *CoRR*, vol. abs/2406.11736, 2024.
 - [151] Y. Qu, T. Zhang, N. Garg, and A. Kumar, “Recursive Introspection: Teaching Language Model Agents How to Self-Improve,” 2024. [Online]. Available: <https://arxiv.org/abs/2407.18219>
 - [152] Y. Deng, P. Lu, F. Yin, Z. Hu, S. Shen, Q. Gu, J. Zou, K.-W. Chang, and W. Wang, “Enhancing Large Vision Language Models with Self-Training on Image Comprehension,” 2024. [Online]. Available: <https://arxiv.org/abs/2405.19716>
 - [153] J. Pang, P. Wang, K. Li, X. Chen, J. Xu, Z. Zhang, and Y. Yu, “Language Model Self-improvement by Reinforcement Learning Contemplation,” in *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
 - [154] Y. Dubois, C. X. Li, R. Taori, T. Zhang, I. Gulrajani, J. Ba, C. Guestrin, P. Liang, and T. B. Hashimoto, “AlpacaFarm: A Simulation Framework for Methods that Learn from Human Feedback,” in *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023, 2023*.
 - [155] A. Madaan, N. Tandon, P. Gupta, S. Hallinan, L. Gao, S. Wiegrefe, U. Alon, N. Dzir, S. Prabhunoye, Y. Yang, S. Gupta, B. P. Majumder, K. Hermann, S. Welleck, A. Yazdanbakhsh, and P. Clark, “Self-Refine: Iterative Refinement with Self-Feedback,” in *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023, 2023*.
 - [156] N. Miao, Y. W. Teh, and T. Rainforth, “SelfCheck: Using LLMs to Zero-Shot Check Their Own Step-by-Step Reasoning,” in *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
 - [157] Z. Gou, Z. Shao, Y. Gong, Y. Shen, Y. Yang, N. Duan, and W. Chen, “CRITIC: Large Language Models Can Self-Correct with Tool-Interactive Critiquing,” in *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
 - [158] Q. Zhong, L. Ding, J. Liu, B. Du, and D. Tao, “ROSE Doesn’t Do That: Boosting the Safety of Instruction-Tuned Large Language Models with Reverse Prompt Contrastive Decoding,” 2024. [Online]. Available: <https://arxiv.org/abs/2402.11889>
 - [159] Y. Weng, M. Zhu, F. Xia, B. Li, S. He, S. Liu, B. Sun, K. Liu, and J. Zhao, “Large Language Models are Better Reasoners with Self-Verification,” in *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*. Association for Computational Linguistics, 2023, pp. 2550–2575.
 - [160] Y. Xie, K. Kawaguchi, Y. Zhao, X. Zhao, M.-Y. Kan, J. He, and Q. Xie, “Self-Evaluation Guided Beam Search for Reasoning,” 2023. [Online]. Available: <https://arxiv.org/abs/2305.00633>
 - [161] Y. Yao, H. Wu, Q. Xu, and L. Song, “Fine-grained Conversational Decoding via Isotropic and Proximal Search,” 2023. [Online]. Available: <https://arxiv.org/abs/2310.08130>
 - [162] J. Chen, W. Lin, J. Mei, and B. Byrne, “Control-DAG: Constrained Decoding for Non-Autoregressive Directed Acyclic T5 using Weighted Finite State Automata,” 2024. [Online]. Available: <https://arxiv.org/abs/2404.06854>
 - [163] N. Xu, C. Zhou, A. Celikyilmaz, and X. Ma, “Look-back Decoding for Open-Ended Text Generation,” 2023. [Online]. Available: <https://arxiv.org/abs/2305.13477>
 - [164] Y. Yao, H. Wu, Z. Guo, B. Zhou, J. Gao, S. Luo, H. Hou, X. Fu, and L. Song, “Learning From Correctness Without Prompting Makes LLM Efficient Reasoner,” *CoRR*, vol. abs/2403.19094, 2024.
 - [165] T. Anthony, Z. Tian, and D. Barber, “Thinking Fast and Slow with Deep Learning and Tree Search,” 2017. [Online]. Available: <https://arxiv.org/abs/1705.08439>
 - [166] N. Miao, Y. W. Teh, and T. Rainforth, “Selfcheck: Using llms to zero-shot check their own step-by-step reasoning,” *arXiv preprint arXiv:2308.00436*, 2023.
 - [167] S. An, Z. Ma, Z. Lin, N. Zheng, J.-G. Lou, and W. Chen, “Learning from mistakes makes llm better reasoner,” *arXiv preprint arXiv:2310.20689*, 2023.
 - [168] Z. Li, X. Hu, A. Liu, K. Zheng, S. Huang, and H. Xiong, “Refiner: Restructure retrieval content efficiently to advance question-answering capabilities,” *arXiv preprint arXiv:2406.11357*, 2024.
 - [169] J. Wu, M. Feng, S. Zhang, F. Che, Z. Wen, and J. Tao, “Beyond examples: High-level automated reasoning paradigm in in-context learning via mcts,” *arXiv preprint arXiv:2411.18478*, 2024.
 - [170] L. Yang, Z. Yu, T. Zhang, M. Xu, J. E. Gonzalez, B. Cui, and S. Yan, “Supercorrect: Supervising and correcting language models with error-driven insights,” *arXiv preprint arXiv:2410.09008*, 2024.
 - [171] L. Yang, Z. Yu, B. Cui, and M. Wang, “ReasonFlux: Hierarchical LLM Reasoning via Scaling Thought Templates,” 2025. [Online]. Available: <https://arxiv.org/abs/2502.06772>
 - [172] Z. Qi, M. Ma, J. Xu, L. L. Zhang, F. Yang, and M. Yang, “Mutual reasoning makes smaller llms stronger problem-solvers,” *arXiv preprint arXiv:2408.06195*, 2024.
 - [173] D. Zhang, J. Wu, J. Lei, T. Che, J. Li, T. Xie, X. Huang, S. Zhang, M. Pavone, Y. Li et al., “Llama-berry: Pairwise optimization for o1-like olympiad-level mathematical reasoning,” *arXiv preprint arXiv:2410.02884*, 2024.
 - [174] L. Yang, Z. Yu, T. Zhang, S. Cao, M. Xu, W. Zhang, J. E. Gonzalez, and B. Cui, “Buffer of Thoughts: Thought-Augmented Reasoning with Large Language Models,” *arXiv preprint arXiv:2406.04271*, 2024.
 - [175] H. Yao, J. Huang, W. Wu, J. Zhang, Y. Wang, S. Liu, Y. Wang, Y. Song, H. Feng, L. Shen et al., “Mulberry: Empowering mllm with o1-like reasoning and reflection via collective monte carlo tree search,” *arXiv preprint arXiv:2412.18319*, 2024.
 - [176] G. Xu, P. Jin, L. Hao, Y. Song, L. Sun, and L. Yuan, “LLaVA-o1: Let Vision Language Models Reason Step-by-Step,” *arXiv preprint arXiv:2411.10440*, 2024.
 - [177] O. Thawakar, D. Dissanayake, K. More, R. Thawkar, A. Heakl, N. Ahsan, Y. Li, M. Zumri, J. Lahoud, R. M. Anwer et al., “LlamaV-o1: Rethinking Step-by-step Visual Reasoning in LLMs,” *arXiv preprint arXiv:2501.06186*, 2025.
 - [178] K. Xiang, Z. Liu, Z. Jiang, Y. Nie, R. Huang, H. Fan, H. Li, W. Huang, Y. Zeng, J. Han et al., “AtomThink: A Slow Thinking Framework for Multimodal Mathematical Reasoning,” *arXiv preprint arXiv:2411.11930*, 2024.
 - [179] Z. Zhang, A. Zhang, M. Li, and A. Smola, “Automatic chain of thought prompting in large language models. arxiv 2022,” *arXiv preprint arXiv:2210.03493*.
 - [180] W. Chen, X. Ma, X. Wang, and W. W. Cohen, “Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks,” *arXiv preprint arXiv:2211.12588*, 2022.
 - [181] L. Gao, A. Madaan, S. Zhou, U. Alon, P. Liu, Y. Yang, J. Callan, and G. Neubig, “Pal: Program-aided language models,” in *International Conference on Machine Learning*. PMLR, 2023, pp. 10 764–10 799.

- [182] T. Khot, H. Trivedi, M. Finlayson, Y. Fu, K. Richardson, P. Clark, and A. Sabharwal, "Decomposed prompting: A modular approach for solving complex tasks," *arXiv preprint arXiv:2210.02406*, 2022.
- [183] D. Zhou, N. Schärli, L. Hou, J. Wei, N. Scales, X. Wang, D. Schuurmans, C. Cui, O. Bousquet, Q. Le *et al.*, "Least-to-most prompting enables complex reasoning in large language models," *arXiv preprint arXiv:2205.10625*, 2022.
- [184] B. Klieger *et al.*, "g1: Using Llama-3.1 70b on Groq to create o1-like reasoning chains," 2024. [Online]. Available: <https://github.com/bklieger-groq/g1>
- [185] X. Hou, M. Yang, W. Jiao, X. Wang, Z. Tu, and W. X. Zhao, "CoAct: A Global-Local Hierarchy for Autonomous Agent Collaboration," *arXiv preprint arXiv:2406.13381*, 2024.
- [186] M. Shen, G. Zeng, Z. Qi, Z.-W. Hong, Z. Chen, W. Lu, G. Wornell, S. Das, D. Cox, and C. Gan, "Satori: Reinforcement Learning with Chain-of-Action-Thought Enhances LLM Reasoning via Autoregressive Search," *arXiv preprint arXiv:2502.02508*, 2025.
- [187] OpenAI, "Reinforcement fine-tuning," 2024.
- [188] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang *et al.*, "A survey on evaluation of large language models," *ACM Transactions on Intelligent Systems and Technology*, vol. 15, no. 3, pp. 1–45, 2024.
- [189] L. Trung, X. Zhang, Z. Jie, P. Sun, X. Jin, and H. Li, "Reft: Reasoning with reinforced fine-tuning," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2024, pp. 7601–7614.
- [190] interconnects.ai, "Blob reinforcement fin-tuning," 2024.
- [191] K. Gandhi, D. Lee, G. Grand, M. Liu, W. Cheng, A. Sharma, and N. D. Goodman, "Stream of search (sos): Learning to search in language, 2024," URL <https://arxiv.org/abs/2404.03683>, 2024.
- [192] K. Team, A. Du, B. Gao, B. Xing, C. Jiang, C. Chen, C. Li, C. Xiao, C. Du, C. Liao *et al.*, "Kimi k1.5: Scaling Reinforcement Learning with LLMs," *arXiv preprint arXiv:2501.12599*, 2025.
- [193] K. Zhang, Q. Yao, B. Lai, J. Huang, W. Fang, D. Tao, M. Song, and S. Liu, "Reasoning with reinforced functional token tuning," *arXiv preprint arXiv:2502.13389*, 2025.
- [194] Z. Lin, Y. Tang, X. Yao, D. Yin, Z. Hu, Y. Sun, and K.-W. Chang, "QLASS: Boosting Language Agent Inference via Q-Guided Stepwise Search," 2025. [Online]. Available: <https://arxiv.org/abs/2502.02584>
- [195] G. Cui, L. Yuan, Z. Wang, H. Wang, W. Li, B. He, Y. Fan, T. Yu, Q. Xu, W. Chen *et al.*, "Process Reinforcement through Implicit Rewards," *arXiv preprint arXiv:2502.01456*, 2025.
- [196] M. Luo, S. Tan, J. Wong, X. Shi, W. Tang, M. Roongta, C. Cai, J. Luo, T. Zhang, E. Li, R. A. Popa, and I. Stoica, "DeepScaleR: Surpassing O1-Preview with a 1.5B Model by Scaling RL," 2025, notion Blog.
- [197] J. Cheng, L. Li, G. Xiong, J. Shao, and Y. Lv, "Stop gamma decay: Min-form credit assignment is all process reward model needs for reasoning," 2025, notion Blog.
- [198] H. Team, "Open r1: A fully open reproduction of deepseek-r1," <https://github.com/huggingface/open-r1>, 2025, github Project.
- [199] J. Pan, J. Zhang, X. Wang, L. Yuan, H. Peng, and A. Suhr, "TinyZero," <https://github.com/Jiayi-Pan/TinyZero>, 2025, accessed: 2025-01-24.
- [200] Z. Liu, C. Chen, W. Li, T. Pang, C. Du, and M. Lin, "There may not be aha moment in r1-zero-like training — a pilot study," <https://oatllm.notion.site/oat-zero>, 2025, notion Blog.
- [201] Z. Liu, C. Chen, C. Du, W. S. Lee, and M. Lin, "Oat: A research-friendly framework for llm online alignment," <https://github.com/sail-sg/oat>, 2025.
- [202] X. Li, H. Zou, and P. Liu, "Limr: Less is more for rl scaling," *arXiv preprint arXiv:2502.11886*, 2025.
- [203] Z. Xie, L. Chen, W. Mao, J. Xu, L. Kong *et al.*, "Teaching language models to critique via reinforcement learning," *arXiv preprint arXiv:2502.03492*, 2025.
- [204] T. Xie, Z. Gao, Q. Ren, H. Luo, Y. Hong, B. Dai, J. Zhou, K. Qiu, Z. Wu, and C. Luo, "Logic-rl: Unleashing llm reasoning with rule-based reinforcement learning," 2025. [Online]. Available: <https://arxiv.org/abs/2502.14768>
- [205] H. Zhang, J. Yao, C. Ye, W. Xiong, and T. Zhang, "Online-dpo-r1: Unlocking effective reasoning without the ppo overhead," 2025, notion Blog.
- [206] J. Hu, Y. Zhang, Q. Han, D. Jiang, and H.-Y. S. Xiangyu Zhang, "Open-reasoner-zero: An open source approach to scaling reinforcement learning on the base model," <https://github.com/Open-Reasoner-Zero/Open-Reasoner-Zero>, 2025.
- [207] T. Yu, Y. Yao, H. Zhang, T. He, Y. Han, G. Cui, J. Hu, Z. Liu, H.-T. Zheng, M. Sun *et al.*, "RLhf-v: Towards trustworthy llms via behavior alignment from fine-grained correctional human feedback," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 13 807–13 816.
- [208] H. Lee, S. Phatale, H. Mansoor, K. R. Lu, T. Mesnard, J. Ferret, C. Bishop, E. Hall, V. Carbune, and A. Rastogi, "Rlaif: Scaling reinforcement learning from human feedback with ai feedback," 2023.
- [209] Y.-F. Zhang, T. Yu, H. Tian, C. Fu, P. Li, J. Zeng, W. Xie, Y. Shi, H. Zhang, J. Wu *et al.*, "Mm-rlhf: The next step forward in multimodal llm alignment," *arXiv preprint arXiv:2502.10391*, 2025.
- [210] J. Ji, J. Zhou, H. Lou, B. Chen, D. Hong, X. Wang, W. Chen, K. Wang, R. Pan, J. Li, M. Wang, J. Dai, T. Qiu, H. Xu, D. Li, W. Chen, J. Song, B. Zheng, and Y. Yang, "Align anything: Training all-modality models to follow instructions with language feedback," 2024. [Online]. Available: <https://arxiv.org/abs/2412.15838>
- [211] L. Chen, L. Li, H. Zhao, Y. Song, and Vinci, "R1-V: Reinforcing Super Generalization Ability in Vision-Language Models with Less Than \$3," <https://github.com/Deep-Agent/R1-V>, 2025, accessed: 2025-02-02.
- [212] H. Shen, Z. Zhang, Q. Zhang, R. Xu, and T. Zhao, "Vlm-r1: A stable and generalizable r1-style large vision-language model," <https://github.com/om-ai-lab/VLM-R1>, 2025, accessed: 2025-02-15.
- [213] Y. Peng, G. Zhang, X. Geng, and X. Yang, "Lmm-r1," <https://github.com/TideDra/lmm-r1>, 2025, accessed: 2025-02-13.
- [214] X. Wang and P. Peng, "Open-r1-video," <https://github.com/Wang-Xiaodong1899/Open-R1-Video>, 2025.
- [215] Y. Zheng, J. Lu, S. Wang, and Y. Xiong, "EasyR1: An Efficient, Scalable, Multi-Modality RL Training Framework," <https://github.com/hiyouga/EasyR1>, 2025.
- [216] E. Yeo, Y. Tong, M. Niu, G. Neubig, and X. Yue, "Demystifying Long Chain-of-Thought Reasoning in LLMs," *arXiv preprint arXiv:2502.03373*, 2025.
- [217] Z. Hou, P. Du, Y. Niu, Z. Du, A. Zeng, X. Liu, M. Huang, H. Wang, J. Tang, and Y. Dong, "Does RLHF Scale? Exploring the Impacts From Data, Model, and Method," *arXiv preprint arXiv:2412.06000*, 2024.
- [218] J. Kim, D. Wu, J. Lee, and T. Suzuki, "Metastable Dynamics of Chain-of-Thought Reasoning: Provable Benefits of Search, RL and Distillation," 2025. [Online]. Available: <https://arxiv.org/abs/2502.01694>
- [219] Z. Liu, C. Chen, W. Li, T. Pang, C. Du, and M. Lin, "There May Not be Aha Moment in R1-Zero-like Training — A Pilot Study," <https://oatllm.notion.site/oat-zero>, 2025, notion Blog.
- [220] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei *et al.*, "Qwen2.5 technical report," *arXiv preprint arXiv:2412.15115*, 2024.
- [221] H. Luo, L. Shen, H. He, Y. Wang, S. Liu, W. Li, N. Tan, X. Cao, and D. Tao, "O1-Pruner: Length-Harmonizing Fine-Tuning for O1-Like Reasoning Pruning," *arXiv preprint arXiv:2501.12570*, 2025.
- [222] J. Hu, "REINFORCE++: A Simple and Efficient Approach for Aligning Large Language Models," *arXiv preprint arXiv:2501.03262*, 2025.
- [223] J. Muralidharan and T. Thomas, "Deliberate Problem-solving with a Large Language Model as a Brainstorm Aid Using a Checklist for Prompt Generation," *The Journal of the Association of Physicians of India*, vol. 72, no. 5, pp. 89–90, 2024.
- [224] J. Jiang, Z. Chen, Y. Min, J. Chen, X. Cheng, J. Wang, Y. Tang, H. Sun, J. Deng, W. X. Zhao *et al.*, "Technical Report: Enhancing LLM Reasoning with Reward-guided Tree Search," *arXiv preprint arXiv:2411.11694*, 2024.
- [225] F. Lyu *et al.*, "Thinking Claude," 2024. [Online]. Available: <https://github.com/richards199999/Thinking-Claude>
- [226] AI-MO, "Aime 2024," <https://huggingface.co/datasets/AI-MO/aime-validation-aime>, 2024.
- [227] —, "Amc 2023," <https://huggingface.co/datasets/AI-MO/aime-validation-amc>, 2024.
- [228] C. He, R. Luo, Y. Bai, S. Hu, Z. L. Thai, J. Shen, J. Hu, X. Han, Y. Huang, Y. Zhang *et al.*, "Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems," *arXiv preprint arXiv:2402.14008*, 2024.

- [229] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. R. Narasimhan, "SWE-bench: Can Language Models Resolve Real-world Github Issues?" in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=VTF8yNQm66>
- [230] N. Jain, K. Han, A. Gu, W.-D. Li, F. Yan, T. Zhang, S. Wang, A. Solar-Lezama, K. Sen, and I. Stoica, "Livecodebench: Holistic and contamination free evaluation of large language models for code," *arXiv preprint arXiv:2403.07974*, 2024.
- [231] D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R. Bowman, "GPQA: A Graduate-Level Google-Proof Q&A Benchmark," in *First Conference on Language Modeling*, 2024. [Online]. Available: <https://openreview.net/forum?id=Ti67584b98>
- [232] Y. Wang, X. Ma, G. Zhang, Y. Ni, A. Chandra, S. Guo, W. Ren, A. Arulraj, X. He, Z. Jiang *et al.*, "Mmlu-pro: A more robust and challenging multi-task language understanding benchmark," *arXiv preprint arXiv:2406.01574*, 2024.
- [233] S. Yao, H. Chen, J. Yang, and K. Narasimhan, "Webshop: Towards scalable real-world web interaction with grounded language agents," *Advances in Neural Information Processing Systems*, vol. 35, pp. 20 744–20 757, 2022.
- [234] S. Zhou, F. F. Xu, H. Zhu, X. Zhou, R. Lo, A. Sridhar, X. Cheng, Y. Bisk, D. Fried, U. Alon *et al.*, "WebArena: A Realistic Web Environment for Building Autonomous Agents," *arXiv preprint arXiv:2307.13854*, 2023. [Online]. Available: <https://webarena.dev>
- [235] R. Wang, P. Jansen, M.-A. Côté, and P. Ammanabrolu, "ScienceWorld: Is your Agent Smarter than a 5th Grader?" 2022. [Online]. Available: <https://arxiv.org/abs/2203.07540>
- [236] A. Prasad, A. Koller, M. Hartmann, P. Clark, A. Sabharwal, M. Bansal, and T. Khot, "ADaPT: As-Needed Decomposition and Planning with Language Models," in *Findings of the Association for Computational Linguistics: NAACL 2024*, 2024, pp. 4226–4252.
- [237] H. Chen, Z. Fang, Y. Singla, and M. Dredze, "Benchmarking Large Language Models on Answering and Explaining Challenging Medical Questions," *arXiv preprint arXiv:2402.18060*, 2024.
- [238] D. Jin, E. Pan, N. Oufattole, W.-H. Weng, H. Fang, and P. Szolovits, "What disease does this patient have? a large-scale open domain question answering dataset from medical exams," *Applied Sciences*, vol. 11, no. 14, p. 6421, 2021.
- [239] X. Yue, Y. Ni, K. Zhang, T. Zheng, R. Liu, G. Zhang, S. Stevens, D. Jiang, W. Ren, Y. Sun *et al.*, "Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 9556–9567.
- [240] P. Lu, H. Bansal, T. Xia, J. Liu, C. Li, H. Hajishirzi, H. Cheng, K.-W. Chang, M. Galley, and J. Gao, "MathVista: Evaluating Mathematical Reasoning of Foundation Models in Visual Contexts," in *International Conference on Learning Representations (ICLR)*, 2024.
- [241] K. Wang, J. Pan, W. Shi, Z. Lu, M. Zhan, and H. Li, "Measuring multimodal mathematical reasoning with math-vision dataset," *arXiv preprint arXiv:2402.14804*, 2024.
- [242] Z.-Z. Li, M.-L. Zhang, F. Yin, Z.-L. Ji, J.-F. Bai, Z.-R. Pan, F.-H. Zeng, J. Xu, J.-X. Zhang, and C.-L. Liu, "Cmmath: A chinese multi-modal math skill evaluation benchmark for foundation models," *arXiv preprint arXiv:2407.12023*, 2024.
- [243] M.-L. Zhang, F. Yin, and C.-L. Liu, "A Multi-Modal Neural Geometric Solver with Textual Clauses Parsed from Diagram," in *IJCAI*, 2023.
- [244] Z. Xi, Y. Ding, W. Chen, B. Hong, H. Guo, J. Wang, D. Yang, C. Liao, X. Guo, W. He, S. Gao, L. Chen, R. Zheng, Y. Zou, T. Gui, Q. Zhang, X. Qiu, X. Huang, Z. Wu, and Y.-G. Jiang, "AgentGym: Evolving Large Language Model-based Agents across Diverse Environments," 2024.
- [245] R. Zhang, D. Jiang, Y. Zhang, H. Lin, Z. Guo, P. Qiu, A. Zhou, P. Lu, K.-W. Chang, P. Gao *et al.*, "MathVerse: Does Your Multimodal LLM Truly See the Diagrams in Visual Math Problems?" *arXiv preprint arXiv:2403.14624*, 2024.
- [246] P. Lu, S. Mishra, T. Xia, L. Qiu, K.-W. Chang, S.-C. Zhu, O. Tafjord, P. Clark, and A. Kalyan, "Learn to Explain: Multimodal Reasoning via Thought Chains for Science Question Answering," in *The 36th Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- [247] Y. Li, Y. Guo, F. Guerin, and C. Lin, "An open-source data contamination report for large language models," in *Findings of the Association for Computational Linguistics: EMNLP 2024*, 2024, pp. 528–541.
- [248] Claude, "Claude 3.5 Sonnet," June 2024. [Online]. Available: <https://www.anthropic.com/news/claude-3-5-sonnet>
- [249] G. DeepMind, "Gemini 2.0 Pro," October 2024. [Online]. Available: <https://deepmind.google/technologies/gemini/pro/>
- [250] I. Team, "InternLM2 Technical Report," 2024.
- [251] N. Muennighoff, Z. Yang, W. Shi, X. L. Li, L. Fei-Fei, H. Hajishirzi, L. Zettlemoyer, P. Liang, E. Candès, and T. Hashimoto, "s1: Simple test-time scaling," *arXiv preprint arXiv:2501.19393*, 2025.
- [252] OpenAI, "OpenAI o1-mini," September 2024. [Online]. Available: <https://openai.com/index/openai-o1-mini-advancing-cost-efficient-reasoning/>
- [253] Y. Ye, Z. Huang, Y. Xiao, E. Chern, S. Xia, and P. Liu, "LIMO: Less is More for Reasoning," 2025. [Online]. Available: <https://arxiv.org/abs/2502.03387>
- [254] Q. Team, "QVQ: To See the World with Wisdom," December 2024. [Online]. Available: <https://qwenlm.github.io/blog/qvq-72b-preview/>
- [255] K. Zhang, B. Li, P. Zhang, F. Pu, J. A. Cahyono, K. Hu, S. Liu, Y. Zhang, J. Yang, C. Li, and Z. Liu, "LMMS-Eval: Reality Check on the Evaluation of Large Multimodal Models," 2024. [Online]. Available: <https://arxiv.org/abs/2407.12772>
- [256] O. Contributors, "OpenCompass: A Universal Evaluation Platform for Foundation Models," <https://github.com/open-compass/opencompass>, 2023.
- [257] M. Song, Z. Su, X. Qu, J. Zhou, and Y. Cheng, "PRMBench: A Fine-grained and Challenging Benchmark for Process-Level Reward Models," *arXiv preprint arXiv:2501.03124*, 2025. [Online]. Available: <https://arxiv.org/pdf/2501.03124>
- [258] Y. Leviathan, M. Kalman, and Y. Matias, "Fast inference from transformers via speculative decoding," in *International Conference on Machine Learning*, 2023, pp. 19 274–19 286.
- [259] X. Ning, Z. Lin, Z. Zhou, Z. Wang, H. Yang, and Y. Wang, "Skeleton-of-thought: Large language models can do parallel decoding," *Proceedings ENLSP-III*, 2023.
- [260] X. Miao, G. Oliaro, Z. Zhang, X. Cheng, Z. Wang, Z. Zhang, R. Y. Y. Wong, A. Zhu, L. Yang, X. Shi *et al.*, "SpecInfer: Accelerating Generative Large Language Model Serving with Tree-based Speculative Inference and Verification," *arXiv preprint arXiv:2305.09781*, 2023.
- [261] B. Qi, X. Chen, J. Gao, D. Li, J. Liu, L. Wu, and B. Zhou, "Interactive continual learning: Fast and slow thinking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 12 882–12 892.
- [262] Y. Zheng, S. Sun, L. Qiu, D. Ru, C. Jiayang, X. Li, J. Lin, B. Wang, Y. Luo, R. Pan *et al.*, "OpenResearcher: Unleashing AI for Accelerated Scientific Research," *arXiv preprint arXiv:2408.06941*, 2024.
- [263] B. Romera-Paredes, M. Barekatin, A. Novikov, M. Balog, M. P. Kumar, E. Dupont, F. J. Ruiz, J. S. Ellenberg, P. Wang, O. Fawzi *et al.*, "Mathematical discoveries from program search with large language models," *Nature*, vol. 625, no. 7995, pp. 468–475, 2024.
- [264] T. H. Trinh, Y. Wu, Q. V. Le, H. He, and T. Luong, "Solving olympiad geometry without human demonstrations," *Nature*, vol. 625, no. 7995, pp. 476–482, 2024.
- [265] Y. Chervonyi, T. H. Trinh, M. Olšák, X. Yang, H. Nguyen, M. Menegali, J. Jung, V. Verma, Q. V. Le, and T. Luong, "Gold-medalist Performance in Solving Olympiad Geometry with AlphaGeometry2," *arXiv preprint arXiv:2502.03544*, 2025.
- [266] G. Team, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, K. Millican *et al.*, "Gemini: a family of highly capable multimodal models," *arXiv preprint arXiv:2312.11805*, 2023.
- [267] G. Team, P. Georgiev, V. I. Lei, R. Burnell, L. Bai, A. Gulati, G. Tanzer, D. Vincent, Z. Pan, S. Wang *et al.*, "Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context," *arXiv preprint arXiv:2403.05530*, 2024.
- [268] N. Chen, Z. Zheng, N. Wu, L. Shou, M. Gong, Y. Song, D. Zhang, and J. Li, "Breaking language barriers in multilingual mathematical reasoning: Insights and observations," *arXiv preprint arXiv:2310.20246*, 2023.
- [269] Y. Du, Z. Liu, Y. Li, W. X. Zhao, Y. Huo, B. Wang, W. Chen, Z. Liu, Z. Wang, and J.-R. Wen, "Virgo: A preliminary exploration on reproducing o1-like mllm," *arXiv preprint arXiv:2501.01904*, 2025.
- [270] M. Parmar and Y. Govindarajulu, "Challenges in Ensuring AI Safety in DeepSeek-R1 Models: The Shortcomings of Reinforcement Learning Strategies," *arXiv preprint arXiv:2501.17030*, 2025.