CoAT: Chain-of-Associated-Thoughts Framework for Enhancing Large Language Models Reasoning

Jianfeng Pan, Senyou Deng, Shaomang Huang

360 Security Technology Inc. Beijing, China

Emails: panjianfeng@360.cn, dengsenyou@360.cn, huangshaomang@360.cn

Abstract

Research on LLM technologies is rapidly emerging, with most of them employ a 'fast thinking' approach to inference. Most LLMs generate the final result based solely on a single query and LLM's reasoning capabilities. However, with the advent of OpenAI-o1, 'slow thinking' techniques have garnered increasing attention because its process is closer to the human thought process. Inspired by the human ability to constantly associate and replenish knowledge during thinking, we developed the novel Chain-of-Associated-Thoughts (CoAT) framework, which introduces an innovative synergy between the Monte Carlo Tree Search (MCTS) algorithm and a dynamic mechanism for integrating new key information, termed 'associative memory'. By combining the structured exploration capabilities of MCTS with the adaptive learning capacity of associative memory, CoAT significantly expands the LLM search space, enabling our framework to explore diverse reasoning pathways and dynamically update its knowledge base in real-time. This allows the framework to not only revisit and refine earlier inferences but also adaptively incorporate evolving information, ensuring that the final output is both accurate and comprehensive. We validate CoAT's effectiveness across a variety of generative and reasoning tasks. Quantitative experiments show that CoAT achieves over 10% performance improvement on open-source multi-hop reasoning datasets (HotpotQA, MuSiQue) and more than 15% gain on our proprietary CRB dataset.

1 Introduction

Large Language Models (LLMs) have rapidly become a cornerstone in natural language processing, powering applications ranging from conversational agents to complex decision-making systems. Central to their operation is the process of inference, where LLMs generate contents by

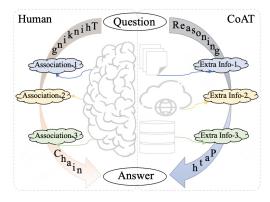


Figure 1: Left: Human thinking chain; Right: Associated thoughts path. This figure illustrates how our CoAT framework is inspired to continually supplement extra information during reasoning by simulating human associative mechanisms.

an auto-regressive learning algorithm in the pretraining stage. Most LLMs, such as GPT (Achiam et al., 2023), LLaMA (Dubey et al., 2024), and Qwen (Yang et al., 2024), employ a 'fast thinking' approach to inference which relies heavily on the pre-trained reasoning capabilities of LLM models. These approaches process a single query to produce the final result. Although effective for many tasks, they often struggle with problems that require nuanced, iterative reasoning, or adaptation to new information.

Recent advances (Li et al., 2022; Brown et al., 2024; Wu et al., 2025) have begun to explore alternatives to 'fast thinking', introducing 'slow thinking' methodologies (Jiang et al., 2024; Min et al., 2024; Gan et al., 2025) that align more closely with human thinking processes. This idea emphasizes deliberate, iterative reasoning, and the integration of historical contents or external knowledge during inference. OpenAI-o1 (Jaech et al., 2024), a notable project, has sparked significant interest in this domain, showcasing the potential of 'slow thinking' frameworks to improve reasoning capabilities. Some studies (Zhang et al., 2023; Liu et al., 2023;

Choi et al., 2023; Chen et al., 2024a; Tian et al., 2024; Zhang et al., 2024; Wan et al., 2024) have employed MCTS-inspired methods to enhance the multi-step reasoning capabilities of LLMs. However, the above mentioned methods merely subdivide the reasoning process into smaller steps and involve rethinking what has already been generated. Throughout the process, reliance is still placed on the initial input information and the logical reasoning abilities of the LLM itself.

Inspired by the human ability to constantly associate and replenish knowledge during thinking, we propose the Chain-of-Associated-Thoughts (CoAT) framework. To our knowledge, associative memory mechanisms were first applied to simulate human thoughts in LLM processes. The associative memory mechanism empowers CoAT to dynamically incorporate new key information during inference, mimicking the human ability to associate and update knowledge iteratively. Furthermore, we optimize the routing strategy in the MCTS algorithm to ensure that each addition of associative memory will provide additional key information for subsequent content generation. This synergy between structured search and adaptive learning enables CoAT to expand its reasoning scope while maintaining contextual coherence, overcoming limitations of conventional LLMs.

The effectiveness of our framework is validated through extensive experiments. The results demonstrate that our framework significantly outperforms traditional models in terms of accuracy, coherence, and diversity. In summary, the main contributions of our work are as follows:

- We propose the CoAT framework to enhance LLM reasoning. Our framework expands the LLM reasoning search space for a better solution using the optimized MCTS algorithm.
- We endow the LLM reasoning process with human-like associative and adaptive selfrefinement capabilities to effectively address complex reasoning tasks.
- We optimize the routing strategy in CoAT to identify the best generation trajectory. The qualitative and quantitative experimental results demonstrate its superior performance compared to other methods.

2 Related Work

The development of Large Language Models (LLMs) has witnessed significant advances in recent years, with a particular focus on improving reasoning capabilities. This section reviews key research on LLM inference strategies, the integration of iterative reasoning frameworks, and associative memory mechanisms, all of which inform the design of our Chain-of-Associated-Thoughts (CoAT).

LLM Inference Strategies Traditional LLMs, including BERT (Devlin, 2018), GPT-3 (Brown et al., 2020) and its successors (like GPT-4 (Achiam et al., 2023)) rely on a single-shot or few-shot inference paradigm. These methods emphasize the model's ability to provide accurate responses using fixed prompts, often resulting in outputs that lack robustness in scenarios that require deeper reasoning. To address these limitations, researchers have explored chain-of-thought (CoT) prompting (Wei et al., 2022) and interleaving retrieval with chainof-thought (IRCoT) (Trivedi et al., 2022a), which enable LLMs to decompose complex problems into smaller sequential steps. Although this improves reasoning quality, it remains inherently static as the model cannot revisit or refine previous inferences during the reasoning process.

More recently, the variants of CoT, such as self-consistency chain-of-thought (CoT-SC) (Wang et al., 2022) have introduced diversity in reasoning by sampling multiple outputs and selecting the most consistent solution, Graph-of-thought (GoT) (Besta et al., 2024) has been improved with search algorithms that can search solution paths more effectively, and Tree-of-thought (ToT) (Yao et al., 2024) prompting uses DFS or BFS search guided by LLMs. However, these methods do not fundamentally alter the underlying inference mechanism, leaving room for further exploration of dynamic and iterative reasoning processes.

The concept of 'slow thinking' (de Winter et al., 2024) has gained traction as an alternative to traditional inference paradigms, inspired by the human ability to deliberate and refine thoughts over time. OpenAI-o1 (Jaech et al., 2024) has been a pioneering framework in this space, demonstrating the benefits of iterative reasoning for tasks involving complex problem solving and decision making. By allowing LLMs to reassess previous steps and integrate new information, slow thinking frameworks improve adaptability and output quality. These advancements highlight the potential

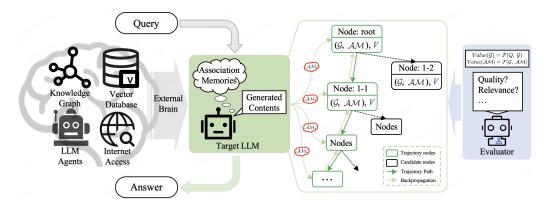


Figure 2: Overview of CoAT framework. The Associative Memory (AM) will be added into each node during reasoning. The "External Brain (\mathbb{EB}) " is an optional measure to further improve the quality of reasoning results.

of moving beyond static reasoning toward more dynamic, context-aware methodologies.

Monte Carlo Tree Search in Inference

MCTS has a long history of success in domains requiring decision making under uncertainty, such as game playing (Silver et al., 2016) and planning (Coulom, 2006). Its ability to balance exploration and exploitation makes it a compelling candidate for enhancing LLM reasoning. Existing works, like LLM-MCTS (Zhao et al., 2024), LLM agent tree search (LATS) (Zhou et al., 2023) and reasoning via planning (RAP) (Hao et al., 2023) have integrated MCTS into specific AI systems to improve search space exploration, but its application in LLMs remains limited. Our CoAT extends this approach by leveraging MCTS not only for structured exploration but also as a means to iteratively refine reasoning pathways by inserting associative memory during inference.

External Knowledge Augmented Mechanisms

Augmented knowledge, an external information retrieval process that enables humans to form and retrieve connections between related concepts when thinking, has inspired various machine learning models. Memory-augmented neural networks (Santoro et al., 2016) and recurrent memory-based architectures (Zaremba, 2014) have demonstrated their effectiveness in tasks requiring long-term context retention. However, these systems often lack the flexibility to adapt to evolving information during LLM inference.

Recent advancements (Gao et al., 2023; Yu et al., 2023; Shao et al., 2023; Chen et al., 2024b; Fan et al., 2024), such as native Retrieval Augmented Generation (NativeRAG) (Lewis et al., 2020), Knowledge Augmented Generation (KAG) (Liang et al., 2024) and hippocampal indexing RAG (Hip-

poRAG) (Gutiérrez et al., 2024), have addressed this by incorporating external knowledge from vector database or knowledge graph at input stage. CoAT framework builds upon this foundation by introducing a dynamic associative memory mechanism that not only retrieves relevant information but also updates and integrates new knowledge in real time during the reasoning stage without requiring post-training. Similarly, Search-R1 (Jin et al., 2025), which is conceptually aligned with CoAT, introduces adaptive retrieval capabilities via the reinforcement learning process, which leads to increased computational costs.

Although existing research has made substantial strides in enhancing the reasoning capabilities and adaptability of LLMs, some challenges remain. Static inference strategies and the limited integration of iterative mechanisms continue to constrain the capacity of LLMs to effectively address increasingly complex and dynamic reasoning tasks. To address these challenges, our proposed CoAT framework synergistically integrates the structured exploration offered by MCTS and the adaptive capabilities of associative memory.

3 Methodology

Inspired by the human ability to form associations during cognitive processes and the demonstrated effectiveness of MCTS algorithm in enhancing the reasoning capability of LLMs, we propose the CoAT reasoning framework, as illustrated in Figure 2. The framework leverages the association mechanism to enable LLMs to perform real-time retrieval of relevant information and self-augmentation during the reasoning process. The realization of this functionality is underpinned by our optimized MCTS algorithm, which systemati-

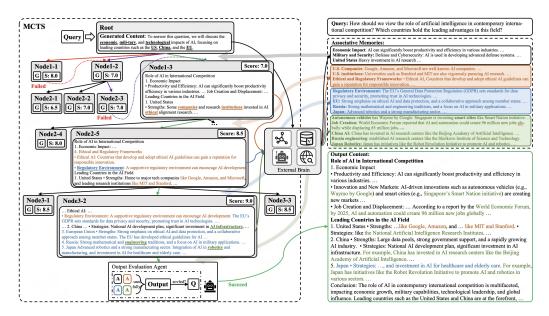


Figure 3: The detailed reasoning process of the CoAT framework. The number of candidate nodes was set to 3.

cally integrates associative content and generated content through tree node search. By assigning precise values to each node based on our predefined rules, the algorithm facilitates the automatic association process, thereby completing the reasoning task. To further enhance the reasoning quality of CoAT framework, we have designed a flexible mechanism for sourcing associative content. This mechanism allows the model to either perform self-association or retrieve associative information through external knowledge sources, referred to as an "External Brain (\mathbb{EB})". The external brain encompasses commonly used resources such as knowledge graph, vector database, LLM agents, and web search engines. A detailed search process of the CoAT framework when query "How should we view the role of artificial intelligence in contemporary international competition? Which countries hold the leading advantages in this field?" is shown in Figure 3.

3.1 Associative Memory Mechanism

We introduce associative memory mechanism in the CoAT framework, can be regarded as a novel external knowledge augmentation mechanism, which enables the reasoning process of LLMs to dynamically update and integrate newly retrieved information in real time according to the generated content of each node. Existing methods primarily focus on incorporating extended knowledge into the reasoning process at its initial stage. However, this approach may lead to incorporation of overly broad knowledge, which introduces two significant draw-

backs: (a) an excess of irrelevant information that compromises inference efficiency, and (b) insufficient inclusion of critical content, ultimately degrading inference quality. In contrast, our proposed real-time association mechanism, integrated into the inference process, effectively addresses these issues by dynamically aligning relevant knowledge with the ongoing inference.

The associative memory mechanism generates content that is beneficial for reasoning and has not been previously mentioned in historical contents. The associative content should exhibit minimal redundancy with existing generated contents and should be concise enough to avoid interfering with the reasoning process. Furthermore, the subject of associative content must maintain a strong relevance to the overall reasoning framework. If these conditions are not satisfied, the associative content for the node can be left empty. The above principle will be applied in evaluation stage for evaluating the quality of associative memory.

When generating the associative memory of a node n_i , the "External Brain" can serve as an alternative approach to enhance the quality of inference results. However, this approach may reduce the efficiency of the inference. This process can be summarized as follows:

$$\mathcal{AM}(n_i) = \mathbb{EB} \mapsto \mathcal{LLM}(Q \mid \mathcal{G}(n_i)). \quad (1)$$

where $\mathcal{G}(n_i)$ denotes the content generated from node n_i and \mathbb{EB} is the External Brain.

Then, a node can reference both the historical content and the associative memories derived from

all of its ancestral nodes. Their historical content and associative content together constitute the comprehensive thinking process of the target LLM. The generation process of each node n_{i+1} is formulated as follows:

$$\mathcal{G}(n_{i+1}) = \mathcal{LLM}(Q \mid \mathcal{G}(n_i) \mid \mathcal{AM}(n_{1:i})). \quad (2)$$

where Q is the input query and $\mathcal{AM}(n_{1:i})$ denotes the associative memories of nodes $n_1 \sim n_i$ in the reasoning trajectory.

3.2 Optimized MCTS

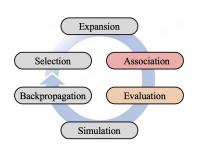


Figure 4: The optimized MCTS process in CoAT.

The standard process of the **MCTS** algorithm consists four stages: Selection, Expansion, Simulation, and Backpropagation. In the selec-

tion stage, MCTS applies the UCT algorithm (Upper Confidence bounds applied to Trees) (Kocsis and Szepesvári, 2006) to choose the best node and then adds it to the trajectory. The UCT of a node n is calculated as follows:

$$UCT(n) = V(n) + w\sqrt{\frac{lnN(p)}{N(n)}}.$$
 (3)

where N(n) is the number of visits to node n, V(n) is the score value, and p is the parent node of node n. w is the exploration weight and is set to 1.0 during CoAT reasoning. When the end of an episode is reached, a back-propagation is carried out to update the value of node n and its parent nodes.

The traditional MCTS algorithm has demonstrated significant success in various decision-making domains. Recently, with advancements in LLM, numerous novel variants of MCTS have been proposed to enable a more effective integration with LLMs. The work of LATS (Zhou et al., 2023) introduces an Evaluation stage after Expansion and a Reflection stage at the end of the process. The evaluation stage assesses the quality of the content generated during the expansion stage, while the reflection stage determines whether the output correctly addresses the inputs. Building on these improvements, we propose an Association stage

to simulate the human associative mechanism between the expansion and evaluation stages. The optimized MCTS process is shown in Figure 4. Consequently, the quality of the associative content is also assessed during the evaluation stage. The evaluation criteria encompass both the quality of the associative content and its correlation with the content generated during the expansion stage, with the goal of preventing excessive associations and mitigating hallucinations. Now, the evaluation value of each node n has two components: the generated content value and the associative content value. And the node value is calculated as follows:

$$V(n) = \mathcal{F}_q(Q, \mathcal{G}(n)) + \beta * \mathcal{F}_a(\mathcal{G}(n), \mathcal{AM}(n)).$$
 (4)

where $\mathcal{G}(n)$, $\mathcal{AM}(n)$ denotes the generated content and the associative content at node n, respectively. \mathcal{F} is the evaluation function for generation and association. β is a weighting coefficient used to control the influence of the associative content, and is set to 0.1 in subsequent experiments.

In the backpropagation stage, we update the visit counts and quality evaluations for every node along the trajectory based on the outcomes of the simulation stage from the leaf node to the root node. The calculation of visit counts is formalized as $C(n_{i+1}) = C(n_i) + 1$. And the quality evaluation value of a parent node n_p will be updated with its children nodes n_c^i as follows:

$$V(n_p)^* = \frac{V(n_p) * C(n_p) + \sum_{i}^{K} V(n_c^i)}{C(n_p) + K}.$$
 (5)

where K is the number of candidate nodes of each parent node, $C(n_p)$ is the original visit counts of n_p . The updated node value $V(n_p)^*$ is used in the UCT algorithm (Eq. 3) to choose the node of the trajectory in the next selection stage.

To more precisely determine when to terminate the MCTS search process, we applied a specialized Reward Model (\mathcal{RM}) to evaluate the content generated at the leaf node of the search trajectory. In certain extreme cases, the search process may enter an ambiguous state, leading to inefficiencies. To mitigate this issue, we introduce a hyper-parameter (D) to constrain the maximum depth of the tree search. When the search depth surpasses D, the process halts, and the best inference result obtained up to that point is returned. Notably, setting D=-1 removes any depth limitation, allowing the search to continue until the optimal result is identified. The flow of the above algorithm can be summarized as Algorithm 1 list at Appendix A.1.

Framework	Model	HotpotQA		2WikiMultiHopQA		MuSiQue	
		EM	F1	EM	F1	EM	F1
NativeRAG	ChatGPT-3.5	43.4 41.8	57.7	33.4	43.3	15.5	26.4
HippoRAG	ChatGPT-3.5		55.0	46.6	59.2	19.2	29.8
IRCoT+NativeRAG	ChatGPT-3.5	45.5	58.4	35.4	45.1	19.1	30.5
IRCoT+HippoRAG	ChatGPT-3.5	45.7	59.2	47.7	62.7	21.9	33.3
IRCoT+HippoRAG	DeepSeek-V2 (236B)	51.0	63.7	48.0	57.1	26.2	36.5
KAG	DeepSeek-V2 (236B)	62.5	76.2	<u>67.8</u>	<u>76.2</u>	36.7	48.7
KAG	Qwen2.5-32B-Instruction	56.6	72.1	65.9	75.5	21.3	31.4
CoAT(Ours)	Qwen2.5-32B-Instruction	69.6	74.2	73.1	78.8	34.7	39.8

Table 1: The end-to-end generation performance of different RAG models on three multi-hop Q&A datasets. The values in **bold** and <u>underline</u> are the best and second best indicators respectively.

4 Experiments

The implementation of our CoAT framework is built upon the LangChain (Contributors, 2025) project. To evaluate the effectiveness of CoAT framework, we designed two types of validation experiments: (a) assessing the qualitative performance of our CoAT framework in conjunction with LLM, via comparative evaluations against baseline models; (b) quantitatively evaluating the CoAT framework against other state-of-the-art reasoning models on both publicly available open-source datasets and custom-constructed complex reasoning benchmarks.

4.1 Qualitative Performance Evaluation

To assess the effectiveness of our CoAT framework in handling real-world reasoning challenges, we designed a series of complex reasoning questions. A case is illustrated in Figure 5. This question requires multidimensional knowledge integration across domains such as economics, ethics. The CoAT-enhanced model (Qwen2.5-32B) outperforms both the baseline Qwen2.5-32B/72B and ChatGPT models, offering more structured and comprehensive responses. Unlike the baseline outputs, which focus on three to four broad categories, the CoAT model organizes its analysis into five clearly defined dimensions: Economic Impact, Military and Security, Technological Leadership, Ethical and Regulatory Frameworks, and Diplomatic and Soft Power. The additional inclusion of the dimension of "Ethical and Regulatory Frameworks" covering AI ethics, privacy regulations, and global governance adds crucial depth and relevance, supported by illustrative examples such as Project Maven and the European AI Alliance.

Moreover, CoAT also demonstrates superior per-

formance in handling the latter part of the task. While baseline models tend to list countries with limited elaboration, CoAT's output delivers a detailed, evidence-based comparison across six items. Each is analyzed in terms of strengths, strategic priorities, and challenges. For instance, the model highlights India's AI initiatives for agriculture and urban development, Japan's roboticcentered AI focus, and Russia's emphasis on military AI within a constrained geopolitical environment. This granular and policy-relevant analysis underscores CoAT's strength in supporting nuanced, multi-perspective reasoning, particularly in domains demanding sophisticated geopolitical insight. The complete response of CoAT can be found in the Appendix A.3.

4.2 Quantitative Performance Evaluation

In this section, we will verify the validity of our CoAT framework in two aspects. (a) We compare the base models' reasoning capacity through the CoAT framework with other retrieval-augmented methods on multi-hop datasets. (b) We compare the results of base models through the CoAT framework with other well-known models on a self-built complex comprehensive reasoning dataset.

Performance on Multi-hop Datasets We enhance the quality of content generated by the associative mechanism through the integration of extended knowledge, and demonstrate that improving the quality of associative content leads to enhanced reasoning ability in our framework. To validate the effectiveness of CoAT framework for the knowledge-intensive question-answering task, we conduct comparative experiments based on retrieval-augmented generation.

The compared methods are NativeRAG (Lewis et al., 2020), IRCoT (Trivedi et al., 2022a), Hip-



Figure 5: The qualitative results of Qwen2.5-32B-Instruct, ChatGPT (APP) and Qwen2.5-72B-Instruct in a complex reasoning query which need rich associative knowledge.

poRAG (Gutiérrez et al., 2024), and KAG (Liang et al., 2024). And three widely-used multi-hop QA datasets are HotpotQA (Groeneveld et al., 2020), 2WikiMultiHopQA (Ho et al., 2020), and MuSiQue (Trivedi et al., 2022b).

Settings. For a fair comparison, we follow IR-CoT, HippoRAG and KAG utilizing a subset of 1,000 questions from each validation set and constructing a retrieval corpus related to selected questions. To evaluate QA performance, we adopt two widely used metrics: Exact Match (EM), and F1 scores. Furthermore, associative memory is influenced not only by the inherent capabilities of the LLM but also by the quality of retrieval results from external knowledge sources. So we apply the results of KAG's retrieval module as the associative memory during CoAT framework evaluation.

Analysis. The multi-hop Q&A performance is presented in Table 1, the results of NativeRAG, HippoRAG and IRCoT using ChatGPT-3.5 and DeepSeek-V2 as the backbone models are excerpts from the official KAG documentation for comparison. However, since the API service for DeepSeek-V2 has been shut down and its local deployment is also costly, we selected the Qwen2.5-32B-Instruction model as an alternative with com-

parable capabilities. Our proposed framework, CoAT, demonstrates significant performance improvements compared to KAG using the same backbone model, with EM gains of 13.0%, 7.2%, and 13.4% on HotpotQA, 2WikiMultiHopQA, and MuSiQue respectively, and F1 improvements of 2.1%, 3.3%, and 8.4%. In particular, the performance of our CoAT with Qwen2.5-32B-Instruction is also better than KAG with DeepSeek-V2, with an increase in EM of 7. 1% and 5. 3% on HotpotQA, 2WikiMultiHopQA datasets.

The observed performance improvements can be largely attributed to the more comprehensive exploration of semantically related entities during the reasoning process in our framework. Leveraging the previous retrieved passages, we employ the association mechanism to identify and expand upon salient entities that are essential for multi-hop reasoning. The retrieval results are further refined in the subsequent content generation stage to enhance response accuracy. However, the association mechanism will lead the model to generate explanatory contents when a direct answer is unavailable, which can reduce response precision and consequently lower the overall F1 score.

Performance on CRB Dataset

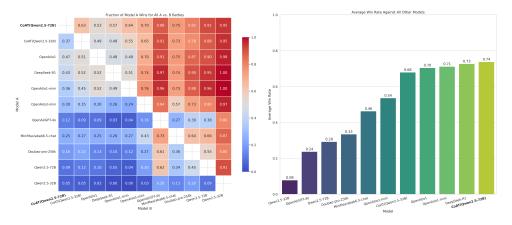


Figure 6: The heatmap of pairwise win rate and the average win rate of all models (Zoom in for best view).

To better demonstrate the effectiveness of CoAT framework in associative reasoning tasks, we constructed a high-quality reasoning dataset, referred to as the Comprehensive Reasoning Benchmark (CRB). This dataset encompasses various disciplines, including politics, scientific and technological domains, international relations, economics, law, and history, among others. The tasks in this dataset require advanced analytical skills, case-based evidence, and rigorous logical reasoning.

Dataset. The CRB dataset contains 205 professionally reviewed questions, each accompanied by its corresponding evaluation rules and total score, which together constitute the final evaluation entries. More details are provided in Appendix A.4 and this dataset has been released in https://github.com/IceSwordLab/CoAT.

Settings. Based on the CRB dataset, we designed two series of experiments. First, we selected multiple state-of-the-art generative models (MiniMax/abab6.5-chat, Doubao-pro-256k, OpenAI/GPT-4o, Qwen2.5-32B-Instruct, Qwen2.5-72B-Instruct), reasoning models (DeepSeek-R1, OpenAI/o1, OpenAI/o1-mini, OpenAI/o3-mini) and our CoAT framework with two base models (Qwen2.5-32B-Instruct, Qwen2.5-72B-Instruct) to generate answers for the questions in the dataset. (For a more detailed discussion on model selection, please refer to the Appendix A.2.) The answers are then evaluated according to Judge Rules, and scores were assigned accordingly. The final average score for each model was computed using a standardized formula: $S_{\mathcal{M}} = \frac{1}{N} \sum_{i}^{\hat{N}} \left(\frac{s_{i}}{s_{T}}\right)$, where s_{i} and s_{T} are the evaluated score and the maximum score, respectively. The experimental results are presented in Table 2.

Second, we conducted pairwise comparisons of

the responses generated by the above models to assess which model's responses demonstrate greater comprehensiveness and depth of detail. Based on these comparisons, we derived the win rate heatmap and the average win rate for each model. The results of this experiment are illustrated in Figure 6.

Models	Average Scores		
Qwen2.5-32B-Instruct	0.55		
OpenAI/GPT-4o	0.59		
Doubao-pro-256k	0.61		
Qwen2.5-72B-Instruct	0.62		
MiniMax/abab6.5-chat	0.66		
OpenAI/o3-mini	0.64		
OpenAI/o1-mini	0.71		
OpenAI/o1	0.73		
DeepSeek-R1	0.75		
CoAT (Qwen2.5-32B)	0.77		
CoAT (Qwen2.5-72B)	0.80		

Table 2: The results of all comparison models on CRB.

Analysis. The results presented in Table 2 support the following conclusions. First, our CoAT framework significantly improves the performance of the base models, Qwen2.5-32B-Instruction and Qwen2.5-72B-Instruction, with relative gains of 22% and 18% in the evaluated scores, respectively. This performance gain can be attributed to two key components: the entity association enabled by the association mechanism, and the progressive exploration strategy of MCTS, which together enhance both the comprehensiveness and accuracy of the generated answers. Second, by leveraging APIs of smaller-scale language models within our CoAT framework, enhanced reasoning performance can be achieved without the need for additional model training or fine-tuning. Finally, reasoning models tend to outperform generative models in CRB datasets, highlighting the advantages of structured reasoning in complex benchmarks.

As illustrated in Figure 6, the Qwen2.5-72B-Instruct model integrated with our CoAT reasoning framework outperforms other models in pairwise evaluations. Specifically, the Owen series achieves a 50% relative improvement in the average win rate when using CoAT, indicating the effectiveness of our proposed framework. Moreover, generative models augmented with the CoAT framework can achieve performance comparable to that of dedicated reasoning models.

4.3 Ablation Experiment

Effectiveness of AM. To separately verify the effectiveness of the associative memory mechanism (Sec. 3.1) and the optimized value computation for MCTS nodes (Sec. 3.2), we performed three experiments using the CRB dataset, and calculated the resulting scores and win rates for both settings: 1) We applied the CoAT framework to generate results either with AM integrated into each node or without it; 2) We incorporated only the content of AM and without considering its quality. 3) We conducted a series of experiments with different values of β in Eq. 4 to explore the optimal solution. In Appendix A.5, we provide an additional detailed discussion of the countermeasures for AM under extreme situations, aiming to enhance the efficiency of the association process.

Models	Versions	Scores
	w/o AM	0.67
CoAT (Qwen2.5-32B)	w/ \mathcal{AM} & β =0	0.75
	w/ \mathcal{AM} & β =0.1	0.77
	w/o $\mathcal{A}\mathcal{M}$	0.71
CoAT (Qwen2.5-72B)	w/ \mathcal{AM} & β =0	0.77
	w/ \mathcal{AM} & β =0.1	0.80

Table 3: The results of the above settings on CRB.

Settings. 1) w/o AM: We disabled the generation of AM at each node as defined in Eq. 1 and set its value to empty in Eq. 2. Subsequently, we omitted the contribution of AM in Eq. 4. 2) w/ \mathcal{AM} & β =**0:** In Eq. 4, β is set to 0.0, while all other components are identical to those in the complete CoAT. The visual results of the above settings are shown in Table 3 and Figure 7. The results of pairwise comparisons between the above settings and all base models are provided in Appendix A.6. 3) The values of β are sampled from [0, 0.6], and the results are shown in Figure 8.

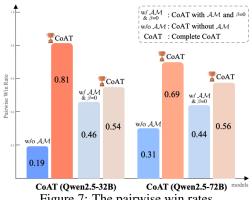


Figure 7: The pairwise win rates.

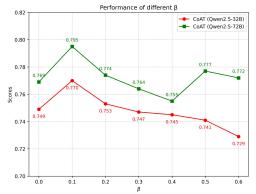


Figure 8: The results of different values of β .

Analysis. The results in Table 3 show that incorporating the whole AM mechanism led to an approximately 10% improvement in average score compared to the baseline without AM. Moreover, the comparison results in Figure 7 further validate the performance improvements achieved through the integration of AM. Figure 8 shows that although the value of β has little impact on the final score, the optimal result is achieved within the entire range when $\beta = 0.1$.

Conclusion 5

In this paper, we proposed the CoAT reasoning framework, which advances LLM reasoning by integrating an optimized MCTS algorithm and a dynamic associative memory mechanism. These innovations enable structured exploration of reasoning pathways and adaptive knowledge updating, addressing limitations of generative LLMs. The experimental results demonstrated that CoAT outperforms other models in accuracy, coherence, and diversity. Our work highlights the potential of combining structured search and adaptive associative memory in LLMs, offering a new exploration for future research on integrating external real-time knowledge for real-world applications.

Limitations

Although our framework outperforms baseline models, there is still room for improvements. Due to the expansion of the search space and the incorporation of novel associative memory, our framework achieves more comprehensive content generation than the baseline models. However, such improvement comes at the cost of increased reasoning time. Additionally, despite the considerable manual efforts invested in curating the self-constructed CRB dataset, there is still room for further quality refinement.

Ethics Statement

Our research focuses on enhancing the reasoning capacity of LLMs. There are no specific ethical concerns directly associated with this work. However, we recognize and emphasize the ethical mindfulness throughout our research. In particular, during the construction of the CRB dataset, no ethical guidelines were violated, and careful attention was paid to data quality and integrity. The broader impact of our work lies in advancing the performance of baseline models, thereby contributing to the improvement of the quality and reliability of content generated by LLMs. All the datasets and models used in this work are publicly available with permissible licenses.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and 1 others. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690.
- Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. 2024. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

- Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. 2024a. Alphamath almost zero: process supervision without process. *arXiv* preprint arXiv:2405.03553.
- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024b. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17754–17762.
- Sehyun Choi, Tianqing Fang, Zhaowei Wang, and Yangqiu Song. 2023. Kcts: knowledge-constrained tree search decoding with token-level hallucination detection. *arXiv* preprint arXiv:2310.09044.
- Contributors. 2025. Langehain: framework for building llm-powered applications.
 - https://github.com/langchain-ai/langchain. Version accessed: langchain==0.3.11, langchain-community==0.3.11, langchain-core==0.3.24.
- Rémi Coulom. 2006. Efficient selectivity and backup operators in monte-carlo tree search. In *International conference on computers and games*, pages 72–83. Springer.
- Joost CF de Winter, Dimitra Dodou, and Yke Bauke Eisma. 2024. System 2 thinking in openai's o1-preview model: Near-perfect performance on a mathematics exam. *Computers*, 13(11):278.
- Jacob Devlin. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. arXiv preprint arXiv:2407.21783.
- Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 6491–6501.
- Zeyu Gan, Yun Liao, and Yong Liu. 2025. Rethinking external slow-thinking: From snowball errors to probability of correct reasoning. *arXiv* preprint *arXiv*:2501.15602.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen Wang, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. arXiv preprint arXiv:2312.10997, 2:1.
- Dirk Groeneveld, Tushar Khot, Ashish Sabharwal, and 1 others. 2020. A simple yet strong pipeline for hotpotqa. *arXiv preprint arXiv:2004.06753*.

- Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. 2024. Hipporag: Neurobiologically inspired long-term memory for large language models. *arXiv preprint arXiv:2405.14831*.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 2023. Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060*.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, and 1 others. 2024. Openai o1 system card. arXiv preprint arXiv:2412.16720.
- Jinhao Jiang, Zhipeng Chen, Yingqian Min, Jie Chen, Xiaoxue Cheng, Jiapeng Wang, Yiru Tang, Haoxiang Sun, Jia Deng, Wayne Xin Zhao, and 1 others. 2024. Technical report: Enhancing Ilm reasoning with reward-guided tree search. *arXiv* preprint *arXiv*:2411.11694.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint* arXiv:2503.09516.
- Levente Kocsis and Csaba Szepesvári. 2006. Bandit based monte-carlo planning. In *European conference on machine learning*, pages 282–293. Springer.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, and 1 others. 2022. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097.
- Lei Liang, Mengshu Sun, Zhengke Gui, Zhongshu Zhu, Zhouyu Jiang, Ling Zhong, Yuan Qu, Peilong Zhao, Zhongpu Bo, Jin Yang, and 1 others. 2024. Kag: Boosting Ilms in professional domains via knowledge augmented generation. *arXiv preprint arXiv:2409.13731*.
- Jiacheng Liu, Andrew Cohen, Ramakanth Pasunuru, Yejin Choi, Hannaneh Hajishirzi, and Asli Celikyilmaz. 2023. Don't throw away your value model! generating more preferable text with value-guided monte-carlo tree search decoding. *arXiv* preprint *arXiv*:2309.15028.

- Yingqian Min, Zhipeng Chen, Jinhao Jiang, Jie Chen, Jia Deng, Yiwen Hu, Yiru Tang, Jiapeng Wang, Xiaoxue Cheng, Huatong Song, and 1 others. 2024. Imitate, explore, and self-improve: A reproduction report on slow-thinking reasoning systems. *arXiv* preprint arXiv:2412.09413.
- Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. 2016. Metalearning with memory-augmented neural networks. In *International conference on machine learning*, pages 1842–1850. PMLR.
- Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy. *arXiv* preprint *arXiv*:2305.15294.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, and 1 others. 2016. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489.
- Qiang Sun, Yuanyi Luo, Wenxiao Zhang, Sirui Li, Jichunyang Li, Kai Niu, Xiangrui Kong, and Wei Liu. 2024. Docs2kg: Unified knowledge graph construction from heterogeneous documents assisted by large language models. *arXiv preprint arXiv:2406.02962*.
- Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Lei Han, Haitao Mi, and Dong Yu. 2024. Toward self-improvement of Ilms via imagination, searching, and criticizing. *Advances in Neural Information Processing Systems*, 37:52723–52748.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022a. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *arXiv* preprint arXiv:2212.10509.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022b. Musique: Multihop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554.
- Ziyu Wan, Xidong Feng, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun Wang. 2024. Alphazero-like tree-search can guide large language model decoding and training. In *Forty-first International Conference on Machine Learning*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *arXiv* preprint arXiv:2203.11171.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances*

- in neural information processing systems, 35:24824–24837.
- Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. 2025. Inference scaling laws: An empirical analysis of compute-optimal inference for llm problem-solving. In *The Thirteenth International Conference on Learning Representations*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.
- Wenhao Yu, Hongming Zhang, Xiaoman Pan, Kaixin Ma, Hongwei Wang, and Dong Yu. 2023. Chain-of-note: Enhancing robustness in retrieval-augmented language models. *arXiv preprint arXiv:2311.09210*.
- Wojciech Zaremba. 2014. Recurrent neural network regularization. *arXiv preprint arXiv:1409.2329*.
- Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. 2024. Rest-mcts*: Llm self-training via process reward guided tree search. *Advances in Neural Information Processing Systems*, 37:64735–64772.
- Shun Zhang, Zhenfang Chen, Yikang Shen, Mingyu Ding, Joshua B Tenenbaum, and Chuang Gan. 2023. Planning with large language models for code generation. *arXiv preprint arXiv:2303.05510*.
- Zirui Zhao, Wee Sun Lee, and David Hsu. 2024. Large language models as commonsense knowledge for large-scale task planning. *Advances in Neural Information Processing Systems*, 36.
- Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. 2023. Language agent tree search unifies reasoning acting and planning in language models. *arXiv preprint arXiv:2310.04406*.

A Appendix

A.1 CoAT Reasoning Algorithm

The flow of our CoAT reasoning algorithm is summarized in Algorithm 1. The details of the CoAT algorithm are described in Section 3.2.

Algorithm 1 CoAT Reasoning Algorithm

Require: Input query (Q), Evaluate function (\mathcal{F}) , Reward Model (\mathcal{RM}) , Number of candidate nodes (K), Search Depth (D), Exploration weight (w), \mathcal{AM} impact factor (β)

```
1: Initialize: Root node n \leftarrow (\mathcal{G}, \mathcal{AM}), d \leftarrow 1
 2: Judge: \mathcal{RM}(\mathcal{G}) is completed for Q.
 3: while Uncompleted and d < D do
        Selection: choose the best trajectory by UCT.
 5:
        for Child Node i \leftarrow 1, ..., K do
            Generation (\mathcal{G}^*): n_i \leftarrow \mathcal{LLM}(Q|\mathcal{G}|\mathcal{AM}).
 6:
            Association (\mathcal{AM}^*): n_i \leftarrow \mathcal{LLM}(Q|\mathcal{G}^*).
 7:
            Evaluation (V): n_i \leftarrow \mathcal{F}(Q, \mathcal{G}^*) + \beta \mathcal{F}(\mathcal{G}^*, \mathcal{AM}^*).
8:
9:
        Judge: \mathcal{RM}(\mathcal{G}^*) of the best child is completed for Q.
10:
        if Completed then
11:
            return \mathcal{G}^*.
12:
        else
13:
            Backpropagation: value of n_p \leftarrow Equation 5.
14:
15:
            Update: visit (V^*) of n_p \leftarrow C(n_p) + K.
            Update: d \leftarrow d + 1.
16:
        end if
17:
18: end while
19: return \hat{\mathcal{G}} of the best trajectory.
```

A.2 Selection of Comparative Models

During the CoAT verification stage, extensive experiments were required, which resulted in a large volume of model access. Considering cost-effectiveness, we prioritized using open-source models. For our task, our primary focus was on the generation and summarization capabilities of LLMs. We researched several mainstream open-source models available at the time. Both the Qwen2.5 series and the Llama3.1 series exhibited comparable performance in these aspects. The detailed comparison results are available on the official Qwen2.5 Blog¹. However, the Qwen series offers broader multilingual support, which is particularly beneficial for the associative memory module, as it may require the understanding of materials in various languages. According to their official Hugging Face model cards, Qwen2.5 supports more than 29 languages², whereas Llama 3.1 supports 8 languages³. Therefore, the Qwen series was deemed more appropriate for our experiments. However, other LLMs (such as the Llama series) can be seamlessly integrated into the CoAT framework simply by replacing the service API or endpoint.

¹https://qwenlm.github.io/blog/qwen2.5

²https://huggingface.co/Qwen/Qwen2.5-72B

³https://huggingface.co/meta-llama/Llama-3.1-70B

A.3 Qualitative Performance Exhibition

The input of all LLMs and CoAT are only the query questions with no extra instruction prompts and system prompts. The responses of three comparison models for a query (How should we view the role of artificial intelligence in contemporary international competition? Which countries hold the leading advantages in this field?) are shown in Figure 5. And the corresponding response of CoAT (Qwen2.5-32B-Instruct) is shown as follows:

```
Role of AI in International Competition
1. **Economic Impact**
    - **Productivity and Efficiency**: AI can significantly boost productivity and
        \hookrightarrowefficiency in various industries. For instance, in manufacturing, AI can
        \hookrightarrowoptimize supply chains and reduce waste. In healthcare, AI-powered
        \hookrightarrow diagnostic tools can improve patient outcomes and reduce costs.
    - **Innovation and New Markets**: AI-driven innovations such as autonomous
         \hookrightarrowvehicles (e.g., Waymo by Google) and smart cities (e.g., Singapore's Smart
        \hookrightarrow Nation initiative) are creating new markets and business opportunities.
    - **Job Creation and Displacement**: While AI can create new jobs in technology
        \hookrightarrowand data science, it can also displace traditional jobs. According to a
        \hookrightarrowreport by the World Economic Forum, by 2025, AI and automation could
        ←create 96 million new jobs globally while displacing 85 million jobs.
2. **Military and Security**
    - **Defense and Cybersecurity**: AI is used in developing advanced defense
        \hookrightarrowsystems. For example, the U.S. Department of Defense has projects like
        \hookrightarrowProject Maven, which uses AI for image recognition in drone footage. China
        \hookrightarrow has also invested in AI for military applications, including autonomous

→drones and cyber defense systems.

    - **Intelligence and Surveillance**: AI can improve intelligence gathering and 

→ surveillance. The U.S. uses AI for analyzing satellite imagery and 

→ detecting potential threats. China has deployed AI in surveillance systems

→, such as facial recognition technology in public spaces.

    - **Strategic Decision-Making**: AI can assist in strategic decision-making. For
        \hookrightarrow example, the U.S. military uses AI to analyze data and provide predictive
        \hookrightarrow insights for military operations.
3. **Technological Leadership**
    - **Research and Development**: Leading countries invest heavily in AI R&D. For
        \hookrightarrowexample, the U.S. has a robust AI research ecosystem with significant
        ←funding from both the government and private sector. China has a national
        →AI development plan with substantial government support.
    - **Patent Filings and Publications**: Leading countries often have a high
        →number of AI-related patent filings and academic publications. According
→to the World Intellectual Property Organization (WIPO), China and the U.S.
        \hookrightarrow are the top two countries in AI patent filings.
    - **Talent Attraction and Retention**: Attracting and retaining top AI talent is
        \hookrightarrow crucial. Countries like the U.S. and Canada offer attractive research
        ←environments and job opportunities. For example, Canada's Vector Institute
        \hookrightarrow is a leading AI research center that attracts top talent from around the
        \hookrightarrowworld.
4. **Ethical and Regulatory Frameworks**
     - **Ethical AI**: Countries that develop and adopt ethical AI guidelines can
        \hookrightarrowgain a reputation for responsible innovation. For example, the EU has
        \hookrightarrow developed ethical guidelines for AI, emphasizing issues like transparency,
        → fairness, and accountability.
    - **Regulatory Environment**: A supportive regulatory environment can encourage
        →AI development. The EU's General Data Protection Regulation (GDPR) sets
        \hookrightarrowstandards for data privacy and security, promoting trust in AI
        \hookrightarrow technologies.
5. **Diplomatic and Soft Power**
     - **Global Standards and Norms**: Leading countries can influence the
        \hookrightarrowdevelopment of global standards and norms for AI. The EU, for instance, is
        \hookrightarrow working on shaping international AI regulations through initiatives like

→ the European AI Alliance.

    - **AI Diplomacy**: Countries can use AI to enhance their diplomatic efforts.
        \hookrightarrow For example, Japan is using AI to improve communication and data-driven
        \hookrightarrow \texttt{decision-making in international negotiations.}
```

Leading Countries in the AI Field

1. United States

- **Strengths**: Home to major tech companies like Google, Amazon, and Microsoft \hookrightarrow , and leading research institutions like MIT and Stanford.

2. China

- **Strategies**: National AI development plan, significant investment in AI →infrastructure, and fostering a large AI talent pool. For example, China →has invested in AI research centers like the Beijing Academy of Artificial → Intelligence.
- **Challenges**: Ensuring data privacy and ethical use of AI, and addressing →international concerns about surveillance and human rights. China faces →scrutiny over its use of AI for social control and surveillance.

3. European Union

- **Strengths**: Strong emphasis on ethical AI and data protection, and a →collaborative approach among member states. The EU has developed ethical →guidelines for AI and promotes data privacy through regulations like GDPR.
- **Strategies**: Development of the General Data Protection Regulation (GDPR)

 → and the European AI Alliance. For example, the EU has launched the

 → European AI Alliance to promote ethical AI development.
- **Challenges**: Harmonizing AI policies across member states and competing →with the US and China in terms of investment and innovation. The EU faces →challenges in coordinating AI policies among its diverse member states.

4. Russia

- **Strengths**: Strong mathematical and engineering traditions, and a focus on →AI in military applications. Russia has a strong research base in →mathematics and computer science.

5. Japan

- **Strengths**: Advanced robotics and a strong manufacturing sector. Japan is a \hookrightarrow leader in robotics and has a strong manufacturing base.
- **Strategies **: Integration of AI in robotics and manufacturing, and →investment in AI for healthcare and elderly care. For example, Japan has →initiatives like the Robot Revolution Initiative to promote AI and →robotics in various sectors.

6. India

- **Strategies **: National AI strategy focusing on healthcare, agriculture, →education, and smart cities. For example, India has launched initiatives →like the National Artificial Intelligence Program to promote AI →development.
- **Challenges **: Infrastructure gaps and the need for significant investment in → AI research and development. India faces challenges in building a robust

The detailed introduction and analysis of the results are described in Section 4.1.

A.4 Introduction of CRB Dataset

The data entry in the CRB⁴ is structured referring to the design principles of subjective questions in the Chinese Gaokao examination. Specifically, each entry consists of three components: Question, Judge Rules, and Score. Judge Rules outline a series of fundamental key points that must be addressed to provide an adequate response. Each key point corresponds to a specific score, and the inclusion of these key points in an answer results in the allocation of the corresponding score. Additionally, Judge Rules incorporate higher-level criteria as bonus points. Score assigned to each data entry represents the maximum attainable score for that entry. Ultimately, we selected 205 professionally reviewed entries as the final test dataset. An example entry from the CRB dataset is shown below:

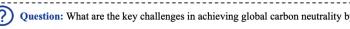
```
question": "What are the key challenges in achieving global carbon neutrality"

ightharpoonupby 2050? Which countries are making the most progress?",
    "judge_rules": "Evaluation Criteria:\n1. (This section is worth 10 points)
        \hookrightarrowDiscuss the main challenges of carbon neutrality, including technological
       \hookrightarrowadditional point can be earned for each aspect if it includes at least one
       → real-world example.\n2. (This section is worth 10 points) Discuss global
       \hookrightarrowleading countries in carbon neutrality efforts, including the European
       \hookrightarrowUnion, China, the United States, Japan, and India. One point is awarded
       →for each country discussed.\n
                                        a. An additional point can be earned for
       \hookrightarroweach country if the discussion covers advantages, disadvantages, and
       ⇒strategies.\n b. Another point can be earned if the discussion of each
       ⇒country includes at least one real-world example.\n3. (This section is
       \hookrightarrowworth 3 points) One point will be awarded for each of the following:
       →fluent language, detailed discussion, and factual accuracy.",
    "total_score": 23
}
```

Each entry in the CRB dataset consists of three components: Question, Judge Rules and Total Score. The language of Question and Judge Rules is either English or Chinese. The key scoring points of the sample shown above are illustrated in Figure 9. For experiments with the CRB dataset, please refer to Section 4.2.

In Figure 9, we provide a detailed breakdown of all the scoring points specified in the judge rules for the question. Each judge rule in the CRB dataset has undergone a professional manual review to assess the reasonableness of the scoring points and the accuracy of the total score. However, there may still be room for further improvement.

⁴https://github.com/IceSwordLab/CoAT



Question: What are the key challenges in achieving global carbon neutrality by 2050? Which countries are making the most progress?



- 1. Key Challenges in Achieving Global Carbon Neutrality by 2050 (10 points)
 - Technological Bottlenecks: 1 point for identification; +1 for real-world example
 - Economic Burden: 1 point for identification; +1 for real-world example
 - Policy Coordination: 1 point for identification; +1 for real-world example
 - Energy Transition: 1 point for identification; +1 for real-world example
 - Social Acceptance: 1 point for identification; +1 for real-world example
- 2. Countries Leading Carbon Neutrality Efforts (10 points)
 - $\ European \ Union: General \ discussion, +1 \ for \ strategy/advantages/disadvantages; +1 \ for \ real-world \ example$
 - China: General discussion, +1 for strategy/advantages/disadvantages; +1 for real-world example
 - United States: General discussion, +1 for strategy/advantages/disadvantages; +1 for real-world example
 - Japan: General discussion, +1 for strategy/advantages/disadvantages; +1 for real-world example
 - India: General discussion, +1 for strategy/advantages/disadvantages; +1 for real-world example
- 3. General Evaluation Criteria (3 points)
 - Fluent language (1 pt)
 - Detailed discussion (1 pt)
 - Factual accuracy (1 pt)



Total Score: 23 points

Figure 9: This figure presents the key scoring points identified in the sample question from the Comprehensive Reasoning Benchmark (CRB) dataset.

A.5 Efficiency Analysis of AM

In our CoAT framework, the association mechanism conducts necessary associations of key entity information based on the generated content of each node to dynamically improve the quality of the final results. Therefore, associating the memory from the initial stage will only involve important information for the key entity information and will not be particularly cumbersome. However, if there is extremely long content in the form of a long document, we will also perform truncation or segmentation by extracting key information, document summaries, or document chunks to ensure that the associative memory actually used for CoAT inference does not become excessively large. Meanwhile, long documents can be transformed into knowledge graphs using document-to-knowledge graph techniques, such as Docs2KG (Sun et al., 2024). Knowledge graphs are inherently supported in CoAT. Moreover, handling long documents is primarily an engineering implementation. The strategies discussed above are capable of addressing this challenge. Therefore, the presence of long texts does not hinder the use of MCTS in CoAT for enhanced reasoning.

Additionally, the following strategies are applicable: (1) Implementing a hierarchical retrieval mechanism by splitting documents into smaller chunks, which allows the EB to return only the most relevant portions for associative reasoning—an approach widely adopted in frameworks like LlamaIndex; and (2) Applying lightweight summarization models such as DistilBART to compress long documents before retrieval, as exemplified in OpenAI's official guide on "Retrieval Augmented Generation (RAG) and Semantic Search for GPTs.⁵"

A.6 Comparison Results of CoAT with Different Settings

The results of the pairwise comparison of the CoAT framework with different settings among all models are shown in Figure 10. The description of this experiment refers to Section 4.3.

⁵https://help.openai.com/en/articles/8868588-retrieval-augmented-generation-rag-and-semantic-sea rch-for-gpts

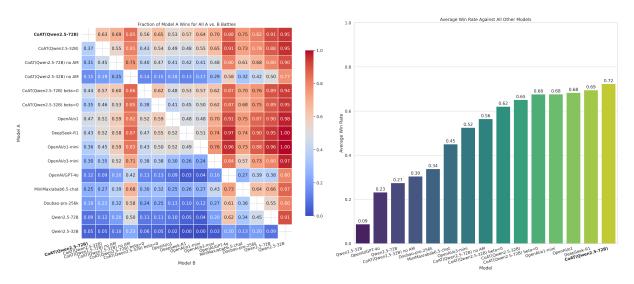


Figure 10: This figure presents the pairwise comparison results among all models. (Zoom in for best view)