

ARM: DISCOVERING AGENTIC REASONING MODULES FOR GENERALIZABLE MULTI-AGENT SYSTEMS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Model (LLM)-powered Multi-agent systems (MAS) have achieved state-of-the-art results on various complex reasoning tasks. Recent works have proposed techniques to automate the design of MASes, eliminating the need for manual engineering. However, these techniques perform poorly, often achieving similar or inferior performance to simple baselines. Furthermore, they require computationally expensive re-discovery of architectures for each new task domain and expensive data annotation on domains without existing labeled validation sets. A critical insight is that simple Chain of Thought (CoT) reasoning often performs competitively with these complex systems, suggesting that the fundamental reasoning unit of MASes, CoT, warrants further investigation. To this end, we present a new paradigm for automatic MAS design that pivots the focus to optimizing CoT reasoning. We introduce the **Agentic Reasoning Module (ARM)**, an agentic generalization of CoT where each granular reasoning step is executed by a specialized reasoning module. This module is discovered through a tree search over the code space, starting from a simple CoT module and evolved using mutations informed by reflection on execution traces. The resulting ARM acts as a versatile reasoning building block which can be utilized as a direct recursive loop or as a subroutine in a learned meta-orchestrator. Our approach significantly outperforms both manually designed MASes and state-of-the-art automatic MAS design methods. Crucially, MASes built with ARM exhibit superb generalization, maintaining high performance across different foundation models and task domains without further optimization.

1 INTRODUCTION

Chain-of-thought (CoT) prompting has emerged as one of the most effective techniques for eliciting complex reasoning from Large Language Models (LLMs) (Wei et al., 2022). By instructing models to generate a series of intermediate steps that lead to a final answer, CoT significantly enhances performance on tasks requiring arithmetic, commonsense, and symbolic reasoning (Nye et al., 2021; Kojima et al., 2022). This simple yet powerful method allows LLMs to break down complex problems into more manageable sub-problems, effectively externalizing the reasoning process over a sequence of generated tokens before arriving at a solution (Wei et al., 2022; Yao et al., 2023a). Recent advancements have also extended CoT with formal verification and multi-agent perspectives, such as MA-LoT (Wang et al., 2025).

Building on the capabilities of individual LLMs, Multi-Agent Systems (MAS) have recently achieved state-of-the-art results on complex reasoning benchmarks (Park et al., 2023; Qian et al., 2023; Hong et al., 2023). These systems typically consist of multiple LLM-powered agents, each assigned a specific role or expertise, orchestrated by a meta-agent or a predefined communication protocol (Wu et al., 2023; Li et al., 2023). While the collaborative nature of MAS enables division of labor and synthesis of diverse perspectives (Chen et al., 2023; Dong et al., 2023), recent work has shifted toward the automatic construction of such systems. Emerging automatic MAS generation frameworks demonstrate how agent roles, communication protocols, and workflows can be synthesized directly by LLMs without manual design. For instance, FlowReasoner and AFlow illustrate this trend by automatically generating agent roles and workflows for LLM-based systems, reducing the need for manual design (Zhang et al., 2025c; Kim et al., 2024).

054 Although MAS approaches have consistently pushed the boundaries of performance, recent studies
055 have revealed a surprising trend: in many cases, a well-prompted single-agent CoT baseline can
056 outperform or perform on par with these complex, multi-agent architectures (Wang et al., 2024; Yao
057 & Yadav, 2025). We also show these observations in our results (Table (1) This finding is significant,
058 as CoT is one of the foundational techniques for LLM reasoning. Its continued competitiveness
059 suggests that the core reasoning unit—the individual thought or step—is of paramount importance.
060 Arguably, the majority of recent research efforts have been dedicated to designing more elaborate
061 MAS frameworks, while the fundamental CoT baseline has remained largely unchanged (Creswell
062 et al., 2022; Chen et al., 2024). Our work pivots from this trend to focus on fundamentally reshaping
063 and enhancing the CoT paradigm for the agentic era by redefining the nature of each reasoning step.

064 In this work, we introduce the Agentic Reasoning Module (ARM), a novel sequential reasoning
065 approach where each granular step is executed by a specialized, self-contained reasoning agent.
066 The core motivation is to elevate the "thinking" steps of CoT from simple textual continuation to
067 the execution of a sophisticated, agentic block. This block is not manually designed but is instead
068 automatically discovered through an evolutionary process. Starting with a basic CoT procedure, the
069 module is iteratively mutated and refined based on its performance on a generic validation dataset
070 of reasoning problems, resulting in a robust and versatile reasoning procedure that can be applied
071 recursively at each step of solving a challenging multi-step problem.

072 The prevailing paradigm for MAS design often leads to systems that are highly domain-specific,
073 with individual agents meticulously tuned for particular skills or tasks (Hu et al., 2025; Zhang et al.,
074 2025b). While single-agent systems are generally considered more versatile, they too are often opti-
075 mized for a narrow set of domains (LaMDAgent, 2025; ScribeAgent, 2024). In contrast, our work
076 focuses on enhancing the universally applicable CoT framework. The agentic block within ARM
077 can be optimized on any generic domain, yielding a general-purpose reasoning technique analogous
078 to the original CoT. We demonstrate that this approach not only achieves superior performance but
079 also exhibits greater generalizability. As we show, MAS built with ARM significantly outperform
080 prominent MAS approaches across diverse agentic datasets without domain-specific tuning.

081 Our methodology uses the simple yet powerful CoT as a starting seed for the evolutionary discovery
082 of ARM. A meta-agent orchestrates this process, performing a tree search over the code space of
083 possible reasoning modules. Mutations and evolutions are guided by a reflection mechanism that
084 analyzes execution traces from previous attempts, identifying weaknesses and proposing targeted
085 improvements. Furthermore, this meta-agent discovers global strategies to orchestrate collabora-
086 tions between parallel ARM reasoning traces, effectively creating a high-performance MAS from
087 optimized, homogeneous building blocks. Overall, our work underscores the immense potential of
088 evolving fundamental reasoning methodologies like CoT, presenting a more robust and scalable al-
089 ternative to the development of increasingly complex and fragile heterogeneous MAS systems. Key
090 contributions of our work are as follows:

- 091 • We present the Agentic Reasoning Module (ARM), an evolved and enhanced version of
092 Chain-of-Thought reasoning. We demonstrate that systems built with ARM substantially
093 outperform existing manually designed and automatically discovered multi-agent systems
094 on complex reasoning tasks.
- 095 • We show that ARM is a significantly more generalizable reasoning module. MAS con-
096 structed with ARM maintain high performance across different underlying foundation mod-
097 els and task domains without requiring re-optimization, highlighting its robustness.
- 098 • We provide a rigorous justification and detailed ablations on the validity of our training
099 objective demonstrating the effectiveness of the proposed MAS discovery strategy.

101 2 RELATED WORKS

102
103 **Single-Agent and Multi-Agent Reasoning Systems** The landscape of LLM-based reasoning is
104 broadly divided into single-agent and multi-agent paradigms. Single-agent systems have demon-
105 strated remarkable capabilities by augmenting the core LLM with sophisticated reasoning and action
106 frameworks. A prominent example is the ReAct framework, which interleaves reasoning steps with
107 actions, enabling the agent to interact with external tools like search engines to gather information
and refine its reasoning process (Yao et al., 2023b). Other approaches have focused on enhancing

108 single agents with self-reflection and memory to learn from past mistakes and improve performance
109 iteratively (Shinn et al., 2023; Madaan et al., 2023). While these systems are powerful, their de-
110 velopment has often focused on narrower tasks, such as tool-based search, retrieval, and question
111 answering, rather than general-purpose complex reasoning.

112 In parallel, Multi-Agent Systems (MAS) have emerged as a dominant approach for tackling highly
113 complex problems, often outperforming single-agent counterparts (Park et al., 2023; Qian et al.,
114 2023). Frameworks like AutoGen (Wu et al., 2023), Camel (Li et al., 2023), and MetaGPT (Hong
115 et al., 2023) orchestrate multiple LLM-powered agents, each assigned a specialized role (e.g., pro-
116 grammer, critic, tester). These agents collaborate, debate, and synthesize information to produce
117 solutions for tasks like software development and complex reasoning. A key characteristic of these
118 systems is their heterogeneous nature; each agent is distinct, with a manually engineered role and
119 persona, connected through a predefined and often complex communication topology. In stark con-
120 trast, our ARM-based approach constructs a powerful MAS from homogeneous building blocks.
121 The ARM itself is a self-contained, versatile reasoning module that is applied repeatedly, acting as
122 the fundamental unit of thought for all "agents" in the system, thereby simplifying the design while
123 enhancing generalizability.

124 **The Surprising Efficacy of Simple Reasoning Baselines** Despite the architectural complexity of
125 many state-of-the-art MAS, a critical and recurring observation is the surprising competitiveness of
126 simple reasoning baselines (Dubey, 2023). Foundational techniques like Chain-of-Thought (CoT)
127 (Wei et al., 2022), and simple extensions like Self-Consistency (CoT-SC) which samples multiple
128 reasoning chains and takes a majority vote (Wang et al., 2022), often achieve performance on par
129 with, or even superior to, intricate multi-agent frameworks (Zhang et al., 2025a). This phenomenon
130 is particularly pronounced with the advent of increasingly powerful frontier foundation models (Ke
131 et al., 2025). As these models develop stronger native reasoning abilities, the high-level conceptual
132 guidance provided by a simple CoT prompt is often sufficient to unlock their full potential, ren-
133 dering the overhead of complex agent orchestration less impactful. This suggests that the primary
134 bottleneck is not necessarily the high-level orchestration strategy but the quality and robustness of
135 the fundamental, step-by-step reasoning process. Our work is directly motivated by this insight,
136 positing that evolving the core reasoning operator—the "thought" in the chain—is a more fruitful
137 direction than designing ever-more-complex superstructures around a static, simple CoT unit.

137 **Automated Design of Multi-Agent Systems** Recognizing the significant manual effort required
138 to design effective MAS, recent research has explored automating this process. Approaches like
139 ADAS (Hu et al., 2025), Aflow Zhang et al. (2025b), and Flow-Reasoner (Gao et al., 2025) aim to
140 automatically discover the optimal agent roles and their interaction topology for a given task do-
141 main. However, these techniques suffer from two major drawbacks. First, they are computationally
142 expensive, requiring a costly re-discovery process for each new task domain. Second, the discov-
143 ered systems are often highly specialized and brittle, tuned specifically for the validation data of a
144 single domain. As our results will later show, with the latest generation of foundation models, these
145 automatically discovered systems can be outperformed by simple CoT baselines. Our work diverges
146 from this paradigm. Instead of discovering a complex, domain-specific agent topology, we focus
147 on discovering a single, domain-agnostic reasoning module (ARM). This ARM acts as a universal,
148 high-quality building block that provides superior performance and generalizability without the need
149 for task-specific rediscovery, offering a more scalable and robust path forward for MAS design.

149 **LLM based Prompt Optimizers** Recent research has focused on LLMs as prompt optimizers,
150 leveraging their generative and reasoning capabilities to automatically improve prompts within a
151 fixed workflow Zhou et al. (2023); Yang et al. (2024); Khattab et al. (2024); Guo et al. (2024);
152 Novikov et al. (2025); Fernando et al. (2024). Notably, evolutionary approaches coupled with deep
153 reflection over rollouts, such as in GEPA Agrawal et al. (2025), have been shown to offer signifi-
154 cant advantages in sample efficiency compared to methods that involve updating model weights via
155 Reinforcement Learning.

157 3 METHODOLOGY: DISCOVERING THE AGENTIC REASONING MODULE

158 We introduce the **Agentic Reasoning Module (ARM)**, a self-contained, code-based multi agentic
159 system designed to execute a single, granular step within a complex reasoning process. ARM is
160 conceived as a structured, agentic replacement for a single step in a Chain of Thought (CoT) se-
161 quence Wei et al. (2022). While standard CoT prompts an LLM to generate the next reasoning step

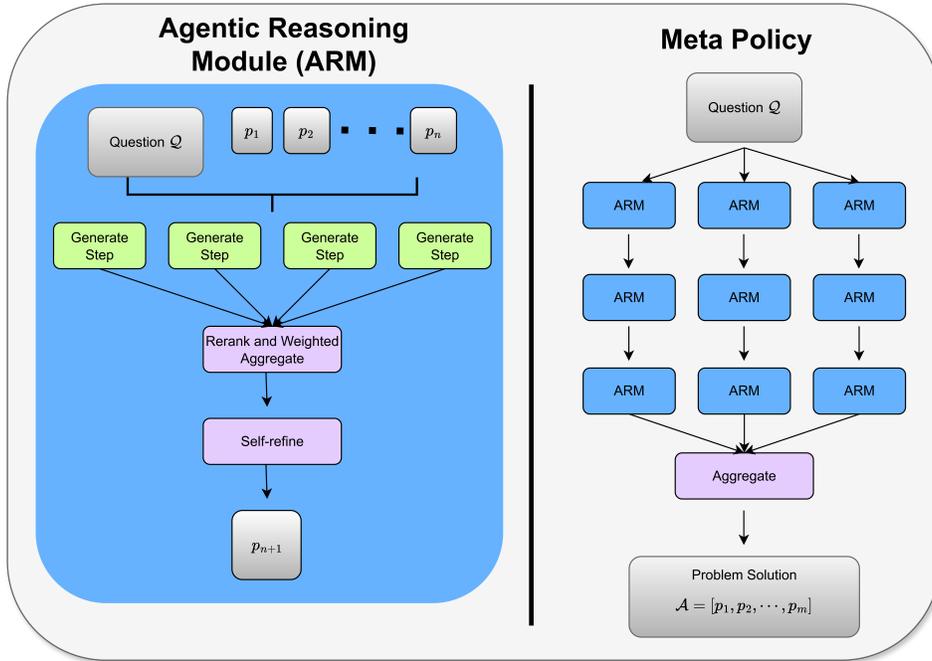


Figure 1: An illustration of the proposed ARM module on the left and the meta policy on the right using "Self refine" as an example MAS. The ARM module takes a question and previous reasoning steps and executes a MAS to get the next step. The meta policy uses ARM as a sub-module and orchestrates the overarching global strategy. Note that this is for illustration only, the actual step generator and the meta policy discovered by Algorithm-1 is more complex (See Appendix).

via naive, monolithic textual generation, an ARM employs an internal multi-agent system (MAS) to produce reasoning steps with greater structure and control.

Following prior work, Hao et al. (2023); Zhang et al. (2024), we define the multi-agentic system as a programming module - a self contained Python function block, while allowing for structured control flow and access to essential APIs such as calling an external LLM, structuring the role and the prompt, and input/output format expectations. Functionally, an ARM accepts the initial problem statement and prior reasoning steps as input, and continues the reasoning until the next logical step in the solution.

3.1 A DECOMPOSABLE FRAMEWORK FOR AGENTIC REASONING

Let the distribution over problem-solution pairs be \mathcal{D} over $(\mathcal{Q}, \mathcal{A})$. A solution \mathcal{A} consists of a sequence of reasoning steps $[p_1, p_2, \dots, p_N]$, where each step p_i belongs to the space of all possible reasoning steps \mathcal{P} . We model the problem-solving process with two key functions:

* **The Step-Generator Module** ($m \in \mathcal{M}$): This is a program that performs a single step of reasoning. It takes the problem question $q \in \mathcal{Q}$ and the history of previous reasoning steps $p_{in} \in \mathcal{P}^*$ as input and returns the next reasoning step $p_{out} \in \mathcal{P}$. Its signature is $m : \mathcal{Q} \times \mathcal{P}^* \rightarrow \mathcal{P}$. An **Agentic Reasoning Module (ARM)** is a structured, code-based implementation of such a module, which can itself be a self-contained MAS.

* **The Meta-Policy** ($\pi \in \Pi$): This is a higher-order program that defines the overarching strategy. It takes a question q and a specific step-generator module m and orchestrates calls to m to generate a complete solution $a \in \mathcal{A}$. Its signature is $\pi : \mathcal{Q} \times \mathcal{M} \rightarrow \mathcal{A}$.

Within this framework, standard Chain of Thought (CoT) can be seen as a simple baseline pairing. It uses a basic step-generator, m_{CoT} , which is a single call to an LLM, and a simple **recursive meta-policy**, π_{Rec} , which applies m_{CoT} repeatedly until a final answer is produced. Our approach independently discovers a more powerful module m^* (the ARM) and a more sophisticated meta-policy π^* .

3.2 DISCOVERING THE OPTIMAL STEP-GENERATOR (m^*)

Our primary goal is to find a step-generator module m^* that is a general-purpose and superior replacement for the simple text generation step in m_{CoT} .

We can formalize a single reasoning step as an update function, $U_{m,q}$, that appends the output of module m to the current reasoning history h :

$$U_{m,q}(h) = h \cdot [m(q, h)]$$

where \cdot denotes list concatenation. A full, n -step reasoning trace generated by the recursive policy π_{Rec} is thus the n -fold composition of this update function: $\pi_{Rec}(q, m) = U_{m,q}^n(\emptyset)$.

Ideally, we would discover the optimal module m^* by maximizing the expected reward \mathcal{R} over the entire problem-solving trace:

$$m^* = \operatorname{argmax}_{m \in \mathcal{M}} \mathbb{E}_{(q,a) \sim \mathcal{D}} [\mathcal{R}(\pi_{Rec}(q, m), a)]$$

However, optimizing this objective directly is intractable due to two main challenges: 1. **Difficult Credit Assignment:** The reward is observed only at the end of a long sequence of steps, making it difficult to determine which specific application of m was responsible for the final outcome. 2. **Unconstrained Search Space:** The space of possible code-based modules \mathcal{M} is vast, making an unguided search highly inefficient.

To address this, we introduce a practical **scaffolded surrogate objective**. Instead of evaluating m on a full rollout generated by itself, we evaluate it within the stable context of a reference trace generated by the baseline m_{CoT} . Specifically, we replace a small, contiguous block of l steps within an n -step CoT trace with our candidate module m . The optimization problem becomes:

$$m^* = \operatorname{argmax}_{m \in \mathcal{M}} \mathbb{E}_{(q,a) \sim \mathcal{D}} [\mathcal{R}(U_{m_{CoT},q}^{n-l-i} \circ U_{m,q}^l \circ U_{m_{CoT},q}^i(\emptyset), a)]$$

where $n = |\pi_{Rec}(q, m_{CoT})|$ is the length of the reference CoT trace, and the starting index i is chosen randomly from $[0, n - 1]$. This formulation isolates the performance contribution of m to a small window, enabling direct credit assignment. Furthermore, the surrounding CoT context provides a powerful inductive bias, constraining the search to modules that behave as effective, incremental reasoning steps. In our experiments, we find $l = 3$ works well, as it is long enough to expose the module m to critical compositional patterns— $(U_{m_{CoT},q} \circ U_{m,q})$, $(U_{m,q} \circ U_{m,q})$, and $(U_{m,q} \circ U_{m_{CoT},q})$ —while keeping the optimization tractable.

3.3 DISCOVERING THE OPTIMAL META-POLICY (π^*)

While an optimized step-generator m^* improves the quality of each reasoning step, the high-level strategy π that orchestrates these steps is equally critical. A simple recursive policy, π_{Rec} , may be suboptimal for complex problems that could benefit from strategies like parallel rollouts (for self-consistency) or iterative refinement loops Wang et al. (2023); Madaan et al. (2023).

Searching for an optimal meta-policy π^* by repeatedly evaluating candidates with the full, complex m^* module is computationally prohibitive. Therefore, we adopt a surrogate-based approach here as well. We search for the optimal meta-policy π^* using the fast and computationally cheap baseline step-generator, m_{CoT} , as a stand-in for m^* .

This zero-shot transfer from m_{CoT} to m^* is effective because our step-generator optimization process (Section 3.2) is explicitly designed to produce an m^* that functions as a superior, "drop-in" replacement for m_{CoT} . A meta-policy that effectively orchestrates the simple steps of m_{CoT} is thus highly likely to generalize to orchestrating the more powerful, but functionally analogous, steps of m^* . This allows us to efficiently explore the space of strategies, discovering sophisticated control flows like branching for parallel thought generation or conditional loops for verification, without incurring the high computational cost of using m^* .

3.4 REFLECTION-GUIDED EVOLUTIONARY SEARCH

We discover both the optimal step-generator m^* and meta-policy π^* using a unified **Reflection-Guided Evolutionary Search** algorithm. This algorithm performs a tree search over the programmatic space of valid Python modules, where each node in the tree represents a specific program.

The search begins with a root node representing the baseline program (m_{CoT} for the step-generator search and π_{Rec} for the meta-policy search). The search then iteratively performs three steps:

1. **Selection:** A parent node (program) p_{parent} is sampled from the current tree \mathcal{T} using temperature sampling based on its validation performance.
2. **Expansion:** A new child program is generated by a **Reviewer Agent**, an LLM-based agent that reflects on the parent program’s execution traces, correctness, and mutation history to propose a targeted code modification.
3. **Evaluation:** The newly generated program is evaluated to obtain its average reward \bar{R} . For a step-generator module, we use the scaffolded objective from Section 3.2. For a meta-policy, we evaluate its performance on a full problem rollout using m_{CoT} as the step-generator.

This entire process is summarized in Algorithm 1.

3.4.1 THE REVIEWER AGENT

The expansion step is driven by a two-stage **Reviewer Agent** that intelligently mutates existing programs. This agent consists of two LLM-based components:

Critic: The Critic analyzes execution traces from the parent program. It identifies logical errors, inefficiencies, or patterns of failure, providing a concise, natural-language analysis of the program’s strengths and weaknesses.

Designer: The Designer acts as the mutation operator. It takes the original program’s code, its performance history, and the Critic’s analysis as input. Based on this information, it proposes a single, targeted code modification aimed at addressing the identified issues, generating a complete, syntactically valid Python class for the new program.

This reflection-driven process ensures that the search evolves programs purposefully, rather than through random mutations, leading to more efficient discovery of high-performance modules and policies. The prompts used for the Critic and Designer are detailed in the Appendix.

4 EXPERIMENTS

4.1 BENCHMARKS

We evaluated our baselines and approach on multiple complex reasoning datasets. To assess complex mathematical reasoning capabilities, we utilized widely studied *American Invitational Mathematics Examination (AIME)*¹ and the *Harvard-MIT Mathematics Tournament (HMMT)*² datasets. For reasoning evaluations on specialized scientific knowledge, we used *GPQA*, a benchmark containing graduate-level questions in physics, chemistry, and biology designed to be challenging even for human experts (Rein et al., 2023). Finally, to measure practical, up-to-date reasoning and robustness against data contamination, we used *LiveBench Reasoning*³, a dynamic benchmark with continuously evolving questions (Kaddour et al., 2023).

4.2 MODELS

We use OpenAI’s o4-mini-high reasoning model as the MAS designer for both the baselines ADAS, AFlow, and our method ARM, as MAS generation requires frontier performance in coding, and instruction following. During validation and inference, we three models as backbone LLMs executing the MAS: two closed source models GPT-4.1-nano, GPT-4o and one open source model Llama-3.3-70B.

4.3 BASELINES

We compare our methodology against two distinct groups of multi-agent systems (MAS) baselines: popular handcrafted MAS systems and leading automated MAS generation approaches.

¹https://huggingface.co/datasets/MathArena/aime_2025

²https://huggingface.co/datasets/MathArena/hmmt_feb_2025

³<https://huggingface.co/datasets/livebench/reasoning>

324 **Handcrafted Multi-Agent Systems:** We compare against several strong reasoning baselines. **Chain**
 325 **of Thought (CoT)** serves as the fundamental baseline, solving tasks through iterative textual rea-
 326 **soning. CoT-Self Consistency (CoT-SC)** improves upon CoT by generating $n = 12$ parallel rea-
 327 **soning rollouts and selecting the final answer via a majority vote. Self-Refine** employs a feedback
 328 **loop where a Large Language Model (LLM) iteratively critiques and refines its own output. Lastly,**
 329 **LLM-Debate** initializes multiple LLM agents with diverse roles to generate different reasoning
 330 **paths, fostering a debate to converge on a final solution.**

331 **Automated Multi-Agent Systems:** These baselines include the two leading code based MAS gen-
 332 **eration approaches: ADAS and AFlow.** These methods employ search algorithms to automatically
 333 **discover the optimal agent roles and their complex interaction topology for a given task domain**
 334

335 Since both ADAS and AFlow are inherently domain-specific and model-specific, we follow the
 336 established methodology by training them on individual validation sets to provide the strongest
 337 possible baseline performance. ADAS runs in an iterative loop with 30 iterations, following the
 338 authors. AFlow runs for until convergence criteria of that method (between 7 to 18 iterations).
 339 We use the AFlow code-base provided by the authors as-is. For a fair comparison on mathematical
 340 questions, we make a minor modification to the ADAS code-base to disallow Python-based symbolic
 341 math tools, instead forcing the model to rely on its inherent reasoning capabilities.

342 4.4 TRAINING

343
 344 Our training process is designed to independently discover the two core components of our frame-
 345 work: the optimal step-generator module (m^*) and the optimal meta-policy (π^*). This decoupled
 346 approach allows us to first forge a powerful, general-purpose reasoning module and then learn a
 347 sophisticated strategy to orchestrate it, all without requiring expensive, domain-specific annotations.
 348

349 **Validation Dataset:** For both discovery processes, we utilize the a subset (1000 samples) of the
 350 Math and Science splits of the Open-or-Mixture-of-Thoughts dataset, a general-purpose instruction-
 351 following dataset. Our method requires only a one-time, domain-agnostic training phase. The same
 352 resulting code artifacts are then deployed across all benchmark domains and foundation models
 353 without any task-specific fine-tuning or re-optimization, underscoring the robustness and versatility
 354 of our method.

355 **Step-Generator (m^*) Discovery:** We discover the ARM module by employing the Reflection-
 356 Guided Evolutionary Search detailed in Algorithm 1. The search is initialized with a basic Chain-of-
 357 Thought module (m_{CoT}) and iteratively evolves it by maximizing the scaffolded surrogate objective
 358 from Section 3.2. This objective evaluates candidate modules within the context of a baseline CoT
 359 trace, enabling efficient and stable optimization.

360 **Meta-Policy (π^*) Discovery:** The meta-policy is discovered independently using the same evolu-
 361 tionary search algorithm. To ensure computational tractability, this search is performed using the
 362 simple and fast baseline module, m_{CoT} , as a surrogate for the more complex m^* (as justified in
 363 Section 3.3). This allows us to efficiently explore the space of high-level strategies and discover a
 364 sophisticated meta-policy that can be seamlessly paired with the optimized ARM module.
 365

366 5 RESULTS

367
 368 We summarize our results in Table 1 and the key findings are as follows:
 369

370 **(1) Naive Operators outperform MAS:** Simple basic operators such as CoT, Self-refine, LLM-
 371 Debate outperform complex MAS systems like AFlow and ADAS. This highlights an important
 372 concern regarding the practicality of recent advancements in MAS. On the other hand, simple
 373 reasoning operators such as CoT perform substantially better across tasks, and varied families of
 374 LLMs. Our ARM based reasoning approach is step forward to revitalize traditional yet strong rea-
 375 soning methods like CoT, by advancing their reasoning steps with agentic blocks. Our ARM based
 376 approach further improves up the CoT performance and achieves best results all the datasets.

377 **(2) ARM achieving top performance:** ARM consistently outperforms all of the operator baselines.
 Specifically, in complex datasets such as AIME and HMMT, ARM consistently outperforms existing

Model	Method	AIME	HMMT	GPQA	LiveBench	Average
GPT-4.1-nano	CoT	15.1%	9.9%	50.0%	33.1%	27.0%
	CoT-SC	<u>21.9%</u>	13.5%	50.6%	36.9%	30.9%
	Self-Refine	17.2%	9.4%	50.0%	28.1%	26.2%
	LLM-Debate	15.1%	<u>16.7%</u>	52.5%	33.8%	29.5%
	ADAS	12.0%	5.2%	48.1%	31.2%	24.1%
	AFlow	18.8%	12.0%	39.9%	30.6%	25.3%
	ARM (Ours)	18.2%	14.6%	<u>60.1%</u>	<u>39.4%</u>	<u>33.1%</u>
	ARM + MP (Ours)	23.4%	22.4%	61.4%	45.6%	38.2%
GPT-4o	CoT	7.3%	0.5%	53.8%	46.2%	27.0%
	CoT-SC	12.5%	2.1%	53.2%	42.5%	27.6%
	Self-Refine	6.8%	2.6%	53.8%	37.5%	25.2%
	LLM-Debate	9.9%	3.1%	56.3%	<u>47.5%</u>	29.2%
	ADAS	1.0%	0.0%	46.2%	38.8%	21.5%
	AFlow	9.9%	3.6%	53.8%	41.9%	27.3%
	ARM (Ours)	<u>13.5%</u>	<u>5.7%</u>	<u>59.5%</u>	<u>47.5%</u>	<u>31.6%</u>
	ARM + MP (Ours)	17.2%	9.4%	60.1%	51.9%	34.7%
LLaMA-3.3-70B	CoT	6.8%	3.1%	50.0%	38.1%	24.5%
	CoT-SC	4.2%	<u>5.7%</u>	53.2%	45.0%	27.0%
	Self-Refine	6.8%	4.2%	<u>51.3%</u>	46.9%	27.3%
	LLM-Debate	5.7%	4.2%	50.6%	46.2%	26.7%
	ADAS	3.1%	0.0%	47.5%	37.5%	22.0%
	AFlow	4.7%	0.0%	46.8%	38.1%	22.4%
	ARM (Ours)	8.3%	5.2%	49.6%	46.2%	<u>27.3%</u>
	ARM + MP (Ours)	<u>7.8%</u>	6.8%	50.0%	50.0%	28.7%

Table 1: Main results on four complex reasoning benchmarks across three foundation models. We compare against two groups of baselines: (1) foundational reasoning strategies used to build agentic systems (CoT, CoT-SC, Self-Refine, and LLM-Debate), and (2) existing state-of-the-art automatic MAS design methods (ADAS and AFlow). Our approach is presented in two variants: **ARM**, which recursively applies the discovered reasoning module, and our full method, **ARM + MP**, which combines the ARM with a learned Meta-Policy (MP). Best score in each category is **bolded** and second best score is underlined.

MAS approaches and all the existing baseline operators. This emphasizes the benefits and strong potential of revitalizing proven traditional reasoning methods like CoT.

(3) Effects from stronger foundation LLM: We first note an important observation that with stronger LLMs such as GPT-4o, simple operators such as CoT and CoT-SC outperform complex MASes. Our ARM based reasoning approach further pushes the best performance over the baseline operators with both recent stronger frontier models such as GPT4.1-nano / GPT-4o and older benchmark models such as LLaMa-3.3-70B.

6 ANALYSES

To understand the sources of ARM’s effectiveness, we performed two key analyses. First, we provide empirical evidence that our search objective discovers fundamentally more reliable reasoning modules by minimizing their per-step error rate. Secondly, we show the validity of our efficient, decoupled training strategy by demonstrating that the learned meta-policy transfers zero-shot from a simple surrogate to the final ARM, yielding significant performance gains.

6.1 EMPIRICAL VALIDATION OF THE STEP-GENERATOR OBJECTIVE

To empirically validate our theoretical claim (Appendix A) that the scaffolded objective minimizes per-step error, we conducted a targeted ablation study. We executed the top five discovered step-generator modules for a single step, starting from *critical reasoning junctures* identified by an LLM-judge (GPT-OSS-20B) within baseline m_{CoT} traces. The error rate of each single-step output was then evaluated. As shown in Figure 1, a module’s rank, determined by our objective, strongly correlates with a lower per-step error rate at these critical points. This result confirms that our search process successfully discovers modules that are fundamentally more robust at a granular level, validating the core mechanism behind ARM’s performance.

6.2 EMPIRICAL VALIDATION OF META-POLICY TRANSFER

Our methodology relies on a crucial transfer: a meta-policy trained with the simple m_{CoT} module is deployed zero-shot with the powerful, discovered m^* module. The theoretical justification in

Meta Policy Name (abbreviated)	CoT Baseline	CoT→Meta	Meta Policy
VWASCCoT	35.1%	33.7%	42.0%
CWDCWACCCoT	37.2%	39.3%	41.8%
RVDCCWASCCoT	33.7%	40.0%	41.8%
DRWASCCoT	35.5%	34.9%	41.8%
MBECDCCWASCCoT	36.3%	39.2%	41.4%

Figure 2: Validation of the meta-policy transfer for top discovered policies. The table compares performance using the simple surrogate m_{CoT} (**CoT Baseline**) versus the powerful ARM module m^* (**Meta Policy**). The intermediate **CoT→Meta** column isolates the performance gain from the superior m^* module alone by evaluating it on states generated by the baseline.

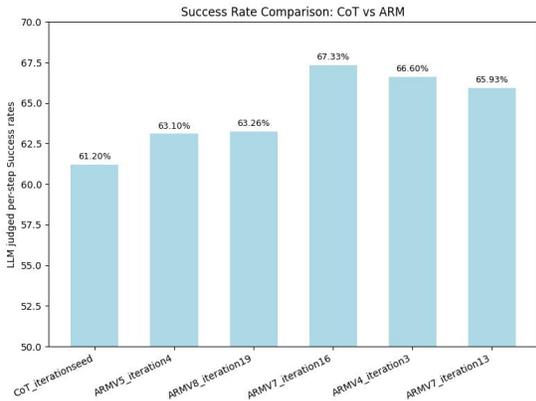


Figure 3: Comparison of LLM judged per-step success rates between the baseline *Chain-of-Thought* (CoT) and multiple *ARM* (CriticChainOfThought) variants. CoT appears first, followed by ARM variants ordered by final performance.

Appendix B posits this transfer is effective due to two factors: (1) the inherent superiority of the m^* module, and (2) its ability to guide the reasoning process into more productive states. We designed an experiment to empirically disentangle and verify these two sources of gain.

To do this, we measure and compare three distinct performance configurations. First, we establish a **baseline performance** using the meta-policy with the simple m_{CoT} module. Second, to isolate the pure **module improvement gain**, we measure the performance of the powerful m^* module when it takes over from intermediate reasoning states generated by the baseline m_{CoT} . Finally, we measure the **full system performance** of the meta-policy paired with m^* from the start.

The results, shown in Figure 2, confirm our hypothesis with a clear performance hierarchy. The baseline system performs worst, followed by a significant improvement from simply swapping to the m^* module. The best performance is achieved by the full system, which benefits from both the better module and its ability to find a better reasoning path. This empirically validates the two conditions for successful transfer outlined in AppendixC and confirms the effectiveness of our decoupled discovery strategy.

7 CONCLUSION

We introduced ARM, a modular agentic reasoning framework that revitalizes the traditional Chain-of-Thought (CoT) paradigm by augmenting it with lightweight agentic blocks. Through extensive experiments, we demonstrated that simple operators such as CoT and Self-Refine not only remain highly competitive but, in many cases, outperform complex Multi-Agent Systems (MAS), highlighting the growing gap between empirical performance and the perceived promise of increasingly elaborate MAS designs. Our results show that ARM consistently advances the performance of CoT across diverse reasoning tasks and model families, establishing top-performing results.

Beyond empirical improvements, ARM sheds light on an important perspective: progress in reasoning with LLMs may not require increasingly complicated MAS architectures, but rather principled extensions of robust, well-understood methods. By preserving the simplicity and generality of the traditional reasoning method—CoT—while enhancing its reasoning depth and modularity, ARM provides a versatile foundation that can be applied across tasks, domains, and models. Overall, this work underscores the importance of revisiting and strengthening proven reasoning approaches instead of overcomplicating them. ARM represents a step toward a practical, scalable, and broadly applicable modular reasoning approach with LLMs, paving the way for future research to place greater emphasis on simple and generalizable traditional methods for complex reasoning.

REFERENCES

- 486
487
488 Lakshya A Agrawal, Shangyin Tan, Dilara Soylu, Noah Ziemis, Rishi Khare, Krista Opsahl-Ong, Ar-
489 nav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, Christopher Potts, Koushik Sen,
490 Alexandros G Dimakis, Ion Stoica, Dan Klein, Matei Zaharia, and Omar Khattab. GEPA: Reflec-
491 tive prompt evolution can outperform reinforcement learning. *arXiv preprint arXiv:2507.19457*,
492 2025.
- 493 Fu-Chieh Chang, Yu-Ting Lee, Hui-Ying Shih, Yi Hsuan Tseng, and Pei-Yuan Wu. RL-star: The-
494 oretical analysis of reinforcement learning frameworks for self-taught reasoner, 2025. URL
495 <https://arxiv.org/abs/2410.23912>.
- 496 Baian Chen, Chang Li, Zhuo Li, Jianing Wang, Yapen Tian, Rui Wang, and Xin Wang. Fireact:
497 Toward language agent fine-tuning. *arXiv preprint arXiv:2403.01925*, 2024.
- 499 Weize Chen, Yusheng Zhang, Zihan Zhang, Cheng Liu, Zipeng Zheng, Chen Qian, Yufan Zhao,
500 Yufan Cong, et al. Agentverse: Facilitating multi-agent collaboration and exploring emergent
501 behaviors. *arXiv preprint arXiv:2308.10848*, 2023.
- 502 Antonia Creswell, Murray Shanahan, and Irina Higgins. Selection-inference: Exploiting large lan-
503 guage models for interpretable logical reasoning. *arXiv preprint arXiv:2205.09712*, 2022.
- 504 Yihong Dong, Xue Wang, Ge Jiang, Zhiping Liu, Cilin Zhang, Peiyu Wang, and Yi Zhang. Self-
505 collaboration code generation via chatgpt. *arXiv preprint arXiv:2304.07590*, 2023.
- 506
507 Akshat Dubey. A simple baseline for multi-agent reasoning. *arXiv preprint arXiv:2310.02023*,
508 2023.
- 509
510 Chrisantha Fernando, Dylan Sunil Banarse, Henryk Michalewski, Simon Osindero, and Tim Rock-
511 täschel. Promptbreeder: Self-referential self-improvement via prompt evolution. In *Proceedings*
512 *of the 41st International Conference on Machine Learning (ICML)*, volume 235, pp. 8370–8386.
513 PMLR, 2024. URL <https://proceedings.mlr.press/v235/fernando24a.html>.
- 514 Hongcheng Gao, Yue Liu, Yufei He, Longxu Dou, Chao Du, Zhijie Deng, Bryan Hooi, Min
515 Lin, and Tianyu Pang. Flowreasoner: Reinforcing query-level meta-agents. *arXiv preprint*
516 *arXiv:2504.15257*, 2025.
- 517
518 Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and
519 Yujiu Yang. Evoprompt: Connecting llms with evolutionary algorithms yields powerful prompt
520 optimizers. In *International Conference on Learning Representations (ICLR)*, May 2024.
- 521
522 Shibo Hao, Yi Gu, Haodi Ma, Joshua Wang, Zhen Chen, and Zhaofeng Wang. Reasoning with
523 language model is planning with world model. *arXiv preprint arXiv:2305.14992*, 2023.
- 524
525 Sirui Hong, Xiawu Zheng, Jonathan Chen, Kechen Yang, Yida Li, Weya Su, Chen Wang, Ceyao
526 He, et al. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint*
arXiv:2308.00352, 2023.
- 527
528 Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems. In *International*
529 *Conference on Learning Representations (ICLR)*, 2025.
- 530
531 Jean Kaddour, Stephen O’Keeffe, Oscar O’Hara, Joe Arts, Sam Sutherland, William Gillingham,
532 Charlie Lee, Leanne d’Avila de Moraes, Adam de Wynter, and Ali Bibi. LiveBench: A real-time
533 llm benchmark. *arXiv preprint arXiv:2311.08277*, 2023.
- 534
535 Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In
536 *Proceedings of the Nineteenth International Conference on Machine Learning, ICML ’02*, pp.
267–274, San Francisco, CA, USA, 2002. Morgan Kaufmann Publishers Inc. ISBN 1558608737.
- 537
538 Zixuan Ke, Fangkai Jiao, Yifei Ming, Xuan-Phi Nguyen, Austin Xu, Do Xuan Long, Minzhi Li,
539 Chengwei Qin, Peifeng Wang, Silvio Savarese, et al. A survey of frontiers in llm reasoning:
Inference scaling, learning to reason, and agentic systems. *arXiv preprint arXiv:2504.09037*,
2025.

- 540 Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vard-
541 hamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller,
542 Matei Zaharia, and Christopher Potts. DSPy: Compiling declarative language model calls into
543 self-improving pipelines. In *The Twelfth International Conference on Learning Representations*
544 *(ICLR)*, 2024.
- 545 Juno Kim, Denny Wu, Jason Lee, and Taiji Suzuki. Metastable dynamics of chain-of-thought rea-
546 soning: Provable benefits of search, rl and distillation, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2502.01694)
547 [abs/2502.01694](https://arxiv.org/abs/2502.01694).
- 548 Sungwoo Kim, Lin Xu, Yifan Guo, Arif Rahman, and Shiyi Wang. Aflow: Automating agentic
549 workflow generation for large language models. In *Proceedings of the 2024 Conference on Em-*
550 *pirical Methods in Natural Language Processing (EMNLP)*, 2024. Accessed: YYYY-MM-DD.
- 551 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
552 language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022.
- 553 LaMDAgent. Lamdagent: An autonomous framework for post-training pipeline construction, 2025.
554 URL <https://arxiv.org/html/2505.21963v1>.
- 555 Guohao Li, Hasan Momin, Hasan Ground, ‘Kian, et al. Camel: Communicative agents for" mind"
556 exploration of large scale language model society. *arXiv preprint arXiv:2303.17760*, 2023.
- 557 Aman Madaan, Niket Tandon, Prakhar Gupta, Kevin Hall, Luyu Gao, Sarah Wiegrefe, Uri Alon,
558 Pengcheng Cair, et al. Self-refine: Iterative refinement with self-feedback. *arXiv preprint*
559 *arXiv:2303.17651*, 2023.
- 560 Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt
561 Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco J. R. Ruiz, Abbas Mehrabian,
562 M. Pawan Kumar, Abigail See, Swarat Chaudhuri, George Holland, Alex Davies, Sebastian
563 Nowozin, Pushmeet Kohli, and Matej Balog. AlphaEvolve: A coding agent for scientific and
564 algorithmic discovery. *arXiv preprint arXiv:2506.13131*, 2025.
- 565 Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Noah Chapman,
566 George Dugan, Miljan Tworkowski, Croitoru Alfredo, et al. Show your work: Scratchpads for
567 intermediate computation with language models. *arXiv preprint arXiv:2112.00114*, 2021.
- 568 Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and
569 Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. *arXiv preprint*
570 *arXiv:2304.03442*, 2023.
- 571 Chen Qian, Xin Wang, Yufan Cong, Cheng Liu, Weize Yu, Zipeng Zheng, Zihan Chen, Yapen Gu,
572 et al. Communicative agents for software development. *arXiv preprint arXiv:2307.07924*, 2023.
- 573 David Rein, Ansh Raichur, Caleb Riddoch, Andrew Andreassen, Ben Jones, Zihui Wu, Shufan
574 Jiang, Kevin Chen, Cong Jiang, Andy Zhao, Lucy Yuan, Jerry Li, Yaofeng Zhang, R Ar-
575 jun Gopalakrishnan, Andrew Pan, Yapei Zhou, Leon Tang, Thomas Lee, Tom Brown, and
576 Jacob Steinhardt. GPQA: A graduate-level google-proof q&a benchmark. *arXiv preprint*
577 *arXiv:2311.12022*, 2023.
- 578 ScribeAgent. Scribeagent: Fine-tuning open-source llms for enhanced web
579 navigation, 2024. URL [https://blog.ml.cmu.edu/2024/12/06/](https://blog.ml.cmu.edu/2024/12/06/scribeagent-fine-tuning-open-source-llms-for-enhanced-web-navigation/)
580 [scribeagent-fine-tuning-open-source-llms-for-enhanced-web-navigation/](https://blog.ml.cmu.edu/2024/12/06/scribeagent-fine-tuning-open-source-llms-for-enhanced-web-navigation/).
- 581 Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: Language agents with verbal rein-
582 forcement learning. *arXiv preprint arXiv:2303.11366*, 2023.
- 583 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford
584 Book, Cambridge, MA, USA, 2018. ISBN 0262039249.
- 585 Qineng Wang, Zihao Wang, Ying Su, Hanghang Tong, and Yangqiu Song. Rethinking the bounds of
586 llm reasoning: Are multi-agent discussions the key? In *arXiv preprint arXiv:2402.18272*, 2024.
587 Accessed: YYYY-MM-DD.

- 594 Ruida Wang, Rui Pan, Yuxin Li, Jipeng Zhang, Yizhen Jia, Shizhe Diao, Renjie Pi, Junjie Hu, and
595 Tong Zhang. Ma-lot: Multi-agent lean-based long chain-of-thought reasoning enhances formal
596 theorem proving. 2025.
- 597 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Self-consistency
598 improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- 600 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha
601 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
602 models. In *International Conference on Learning Representations (ICLR)*, March 2023.
- 603 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
604 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
605 neural information processing systems*, 35:24824–24837, 2022.
- 607 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Li, Erkang Zhu, Beibin Li, Li Jiang,
608 et al. Autogen: Enabling next-gen llm applications via multi-agent conversation. *arXiv preprint
609 arXiv:2308.08155*, 2023.
- 610 Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen.
611 Large language models as optimizers. In *International Conference on Learning Representations
612 (ICLR)*, May 2024.
- 613 Bohan Yao and Vikas Yadav. A toolbox, not a hammer—multi-tag: Scaling math reasoning with
614 multi-tool aggregation. *arXiv preprint arXiv:2507.18973*, 2025.
- 616 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Sha, Silvio Savarese, and Sima an. Tree of thoughts:
617 Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*, 2023a.
- 618 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Sha, Narasimhan Karthik, and Sima an. React:
619 Synergizing reasoning and acting in language models. In *The Eleventh International Conference
620 on Learning Representations*, 2023b.
- 622 Guibin Zhang, Luyang Niu, Junfeng Fang, Kun Wang, Lei Bai, and Xiang Wang. Multi-agent
623 architecture search via agentic supernet. In *Proceedings of the 42nd International Conference on
624 Machine Learning (ICML)*, 2025a. Oral Presentation (Top ~1%).
- 625 Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xiong-Hui Chen, Jiaqi Chen, Mingchen
626 Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, Bingnan Zheng, Bang Liu, Yuyu Luo, and Chenglin
627 Wu. AFlow: Automating agentic workflow generation. In *International Conference on Learning
628 Representations (ICLR)*, 2025b. Oral Presentation (Top 1.8%).
- 629 Lei Zhang, Feng Xu, Zongyi Yu, Chen Zhu, and Yu Qian. Agent-flow: A flexible and scalable
630 multi-agent platform for real-life tasks. *arXiv preprint arXiv:2402.17779*, 2024.
- 632 Wei Zhang, Chen Liu, Ananya Patel, Ming Zhao, and Jie Huang. Flowreasoner: Automatic multi-
633 agent system generation for complex reasoning. *arXiv preprint arXiv:2502.08123*, 2025c. Ac-
634 cessed: YYYY-MM-DD.
- 635 Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and
636 Jimmy Ba. Large language models are human-level prompt engineers. In *International Confer-
637 ence on Learning Representations (ICLR)*, November 2023.
- 638
639
640
641
642
643
644
645
646
647

A ANALYSES

We conducted two practical analyses to evaluate a) stability of ARM over multiple iterations during search process, b) cost comparison to ARM with prior AMAS techniques.

A.1 STABILITY OF ARM DURING SEARCH

We evaluated the performance of ARM versions at each iteration to show convergence and stability of the discovered ARM module. As shown in 4, the performance of the ARM modules in later iterations converges to the final performance.

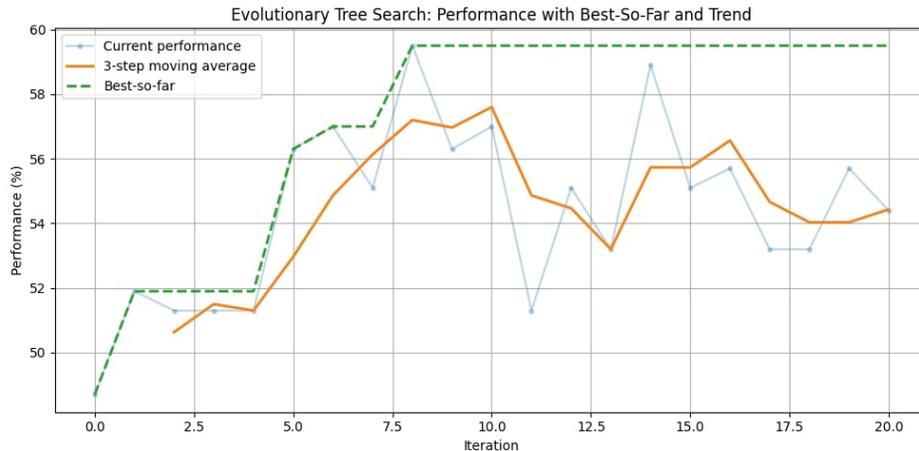


Figure 4: Plot showing performance on GPQA benchmark on intermediate candidates of the search. Note that the search space is discrete (code modules) and the search is tree based exploration/exploitation paradigm, which explains the bumps in performance. However, all intermediate checkpoints are stable, and within 5% of each other.

A.2 SEARCH AND INFERENCE COST ANALYSIS

There are two costs to analyze: **search cost** and **inference cost**.

Model	Cost (USD)
ARM (Ours)	4.53
Meta (Ours)	5.40
Total Cost (Ours)	9.93
ADAS	2.62
AFlow	2.52

Table 2: Training cost per iteration on 1K samples of open-r1 training data using GPT-4.1-mini, our best performing model

A) Training Cost per Iteration on 1K Samples (GPT-4.1-mini) A single iteration of ARM + Meta policy is $\sim 3.8\times$ the cost of ADAS and $\sim 3.9\times$ the cost of AFlow. However, this must be interpreted correctly:

- Domain- and model-specific techniques such as ADAS, AFlow (and MaAS) must be run **independently for each domain and model**.
- Thus, with 4 domains and 3 models, ADAS cost becomes:

$$4 \times 3 \times 2.62 = \$31.44$$

and AFlow cost becomes:

$$4 \times 3 \times 2.52 = \$30.24.$$

- In contrast, ARM is **domain- and model-agnostic**, so its amortized cost stays constant—approximately $0.31 \times$ that of ADAS and $0.32 \times$ that of AFlow.

Cross-Domain Applicability of ADAS and AFlow For fairness, we trained ADAS and AFlow on the same 1K OpenReasoner subset as ARM. Both methods show **poor cross-domain transfer**. Even with model-specific training, we observed no gains when evaluating on unseen domains, indicating that these approaches cannot leverage large-scale datasets beyond the limited benchmark-specific validation sets.

Method	Model	AIME	HMMT	GPQA	LiveBench
ADAS (1K)	GPT-4.1-mini	0.0	6.8	46.8	29.4
ADAS (1K)	GPT-4o	0.0	0.5	46.8	41.9
ADAS (1K)	Llama-3.3-70B	3.1	0.5	42.4	46.2
AFlow (1K)	GPT-4.1-mini	16.7	10.4	51.3	30.6
AFlow (1K)	GPT-4o	9.4	0.0	50.6	45.0
AFlow (1K)	Llama-3.3-70B	7.2	3.1	46.8	15.6

Table 3: Cross-domain performance when trained on 1K OpenReasoner samples.

Each row above represents a **model-specific but domain-agnostic MAS** configuration, in contrast to ARM which is both **domain- and model-agnostic**.

A.2.1 B) INFERENCE COST (USD) ON 1K OPENREASONER SAMPLES (GPT-4.1-NANO)

Method	AIME	HMMT	GPQA	LiveBench
ARM (Ours)	3.69	3.02	0.81	1.12
ARM + Meta (Ours)	17.60	13.77	3.22	4.36
ADAS	0.87	0.88	0.19	0.28
AFlow	0.78	0.66	0.16	0.19

Table 4: Inference cost on 1K samples using GPT-4.1-nano.

ARM is designed as a *compute-for-performance* approach. Instead of running the workflow once—risking severe error propagation—ARM executes the reasoning module at each granular step. This trades compute for significant performance gains and enables automatic discovery of multi-agentic systems.

Cost-Optimized ARM Variants. Since ARM is iteratively refined using a GEPA-style reflection-guided search, adding a cost term to the reviewer agent yields Pareto-optimal cost–performance variants. We will expand on this in the final version of the paper.

Approach	Cost Multiple	Performance Gain	Gain / Cost
CoT	1.0 (\$0.13)	–	–
ADAS	4.2	-2.9	-0.69
AFlow	3.5	-1.7	-0.48
ARM (Ours)	16.6	+6.1	+0.36
ARM + Meta (Ours)	74.8	+11.2	+0.14

Table 5: Cost vs. performance comparison.

Cost–Performance Comparison with CoT and MAS Approaches ADAS, AFlow, and other MAS methods are substantially more expensive than simple CoT yet offer only marginal or negative gains. ARM, while costlier, delivers **significantly higher and consistent performance**, with a more favorable cost-normalized improvement ratio compared to all other automated agentic systems.

B ARM SEARCH ALGORITHM

Algorithm 1 provides the full pseudocode of the reflection-guided search algorithm for evolving ARM modules.

Algorithm 1 Reflection-Guided Search

```

1: Input: Initial program  $p_{root}$  (e.g.,  $m_{CoT}$  or  $\pi_{Rec}$ ), evaluation function  $EVALUATE(\cdot)$ , total
   iterations  $K$ , exploration constant  $C$ .
2: Initialize:
3: Tree  $\mathcal{T}$  with a single node for  $p_{root}$ .
4:  $p_{root}.\bar{\mathcal{R}} \leftarrow EVALUATE(p_{root})$   $\triangleright$  Evaluate the baseline program on a validation batch
5:  $p_{root}.N \leftarrow 1$   $\triangleright$  Initialize visit count for the root
6: for  $t = 1$  to  $K$  do
7:  $\triangleright$  1. Select a parent program to mutate
8:  $P(p_i) \leftarrow \frac{\exp(p_i.\bar{\mathcal{R}}/T)}{\sum_{j \in \mathcal{T}} \exp(p_j.\bar{\mathcal{R}}/T)}$ 
9:  $p_{parent} \leftarrow \text{Sample}(\mathcal{T}, P)$ 
10:  $\triangleright$  2. Expand the tree via reflection
11:  $\text{traces} \leftarrow EXECUTE(p_{parent})$   $\triangleright$  Collect execution traces
12:  $\text{history} \leftarrow \text{GETMUTATIONHISTORY}(p_{parent})$ 
13:  $p_{new} \leftarrow \text{REVIEWERAGENT}(p_{parent}, \text{traces}, \text{history})$ 
14:  $\triangleright$  3. Evaluate the new program
15:  $p_{new}.\bar{\mathcal{R}} \leftarrow EVALUATE(p_{new})$ 
16:  $p_{new}.N \leftarrow 1$ 
17:  $\triangleright$  4. Update tree and statistics
18:  $\mathcal{T}.\text{ADDCHILD}(p_{parent}, p_{new})$ 
19:  $p_{parent}.N \leftarrow p_{parent}.N + 1$ 
20: end for
21:
22: return  $\underset{p_i \in \mathcal{T}}{\text{argmax}} (p_i.\bar{\mathcal{R}})$   $\triangleright$  Return the program with the highest empirical reward

```

C THEORETICAL ANALYSIS

A complete theoretical analysis of the multi-agentic system ARM powered by LLMs is intractable due to the complex, high-dimensional nature of language generation and the non-stationary of the generation process. Recent research (Chang et al., 2025; Kim et al., 2025) models sequential CoT reasoning steps as a Markov Decision Process by abstracting away the underlying complexities of the text generation process and focusing on higher level reasoning states. Therefore, to build a formal intuition for the design choices in our scaffolded search for the step-generator, and the decoupled search for the meta-policy (Algorithm1), we also analyze an idealized formulation of the problem as a Markov Decision Process (MDP).

Our analysis is particularly inspired by recent work on self taught reasoners (RL-Star) by Chang et al. (2025), where they introduce a step indexed competence parameter $\delta_{t,n}$ which quantifies the advantage in probability of a correct reasoning step at step n during training iteration t over a baseline random reasoner. They show the conditions under which a bootstrapped RL learning algorithm based on rejection sampling shows monotonic improvement and convergence. While our goals are similar (improving the reasoning process), our problem statement has critical differences which makes a straight forward adaption infeasible: RL-Star analyses a system where the LLM’s parametric weights are updated via reinforcement learning. On the other hand, ARM treats the LLM as a black box and performs discrete, evolutionary search (Fernando et al., 2024; Agrawal et al., 2025) over programming modules that orchestrate calls to the LLM. Consequently, our search is inherently *discrete*, so smoothness-based guarantees do not apply. Hence we do not assume or prove convergence guarantees, and instead motivate the intuition of our scaffolded search process as a conservative policy improvement (CPI) (Kakade & Langford, 2002) that preferentially selects modules with higher competence leading to improved reasoning process.

C.1 AN IDEALIZED MDP MODEL OF STEP-WISE REASONING

We model the reasoning process as a Markov decision process (MDP) Sutton & Barto (2018) $\mathcal{M} = (S, N, A, P, R, \gamma)$:

- **State Space (S):** The state space $S = \mathcal{U} \cup \mathcal{G} \cup \mathcal{F}$ is partitioned into three disjoint subsets:
 - \mathcal{U} : A state $s \in \mathcal{U}$ represents a partial reasoning trace q, p_1, \dots, p_k that is not yet terminated.
 - \mathcal{G} : A state $s \in \mathcal{G}$ represents a reasoning path that has successfully ended on the right answer. In our setting this is when the module emits the `/boxed{correct answer}`. This is an absorbing region.
 - \mathcal{F} : A state $s \in \mathcal{F}$ represents a reasoning path that has terminated at the wrong answer. In our setting this is when the module emits the `/boxed{incorrect answer}`. This is an absorbing region.
- **Verification Predicate (`solved`):** A predicate function $S \rightarrow \{0, 1\}$ judging if the right answer is already derivable from the given partial reasoning state. Note that this is a simple formatting action, and is independent of the module m .
 - `solved`(s) = 0 $\forall s \in \mathcal{F}$
 - `solved`(s) = 1 $\forall s \in \mathcal{G}$
- **Maximum Reasoning Steps (N):** We rollout the reasoning process up to N steps. After N steps of reasoning, we enforce a *model-independent* termination rule where the state deterministically goes into $s' \in \mathcal{G}$ if `solved`(s) = 1 and into $s' \in \mathcal{F}$ if `solved`(s) = 0. For simplicity of notation, we assume the total trajectory length to be $N + 1$.
- **Action Space (A):** For a fixed meta-policy π_{Rec} that recursively generates steps until termination (such as the one used by baseline CoT or the ARM-only variant), the meta policy executes a single action at any give state $s \in \mathcal{U}$: i.e., invokes a step-generator module m to produce the next reasoning step. Thus, the action space is a singleton $\mathcal{A} = \{\text{generate_step}\}$. For terminal states $\mathcal{G} \cup \mathcal{F}$, this is a `no-op`. Hence, the choice of the module m fully defines the transition dynamics of the MDP.
- **Reward Function (R):** The one-shot terminal reward is sparse:

$$R(s \rightarrow s') = \begin{cases} 1, & s' \in \mathcal{G}, \\ 0, & \text{otherwise.} \end{cases}$$

- **Transition Dynamics (P):** We denote the state transition probability $P(\cdot|s, m)$ with the Markov assumption. This simplification is the core foundation for our MDP analysis.
- **Value Function:** For $n \in \{0, \dots, N\}$, let $V_n^m(s)$ denote the value with n reasoning steps remaining before the formatting step. The Bellman recursion can be written as

$$V_n^m(s) = \begin{cases} 1, & s \in \mathcal{G}, \\ 0, & s \in \mathcal{F}, \\ \{\text{solved}(s)\} + \{\neg \text{solved}(s)\} \mathbb{E}_{s' \sim P_m(\cdot|s)} [V_{n-1}^m(s')], & s \in \mathcal{U}, n \geq 1, \\ \{\text{solved}(s)\}, & s \in \mathcal{U}, n = 0. \end{cases}$$

Within this MDP framework, the ideal objective is to discover a module m^* , that maximizes the expected value from the initial state distribution $d_0(s)$

$$m^* = \arg \max_{m \in \mathcal{M}} \mathbb{E}_{s_0 \sim d_0(s)} [V_N^m(s_0)]$$

This objective poses several major optimization challenges: 1) credit assignment problem over long sequence of steps and 2) unconstrained search space of code modules.

864 C.2 DEFINITIONS

865 We introduce the following quantities to characterize module’s performance and search strategy:

- 866 • **Per-step competence** $\delta_m(s)$: This represents the *competence* of the module m at a reason-
867 ing state $s \in \mathcal{U}$ analogous to $\delta_{t,n}$ term in Chang et al. (2025). The probability that a
868 one step update $s \in \mathcal{U}$ is valid can be viewed as a monotonically increasing function over
869 $\delta_m(\cdot)$, but for simplicity of notation, we assume this to be $\delta_m(s)$ itself.
- 870 • **Recovery** $r_m(s)$: The probability that, conditioned on an invalid one-step update from s ,
871 the next step returns to a valid state. This term captures the recovery possibility of a mistake
872 in the immediate next turn. While recovery can happen at any turn following the mistake
873 in a real LLM, we limit the window to 1 for simplicity.
- 874 • **Composite Validity** $\phi_m(s)$: The total probability that the next step is valid for a step $s \in \mathcal{U}$,
875 either by being immediately valid, or by being successfully repaired on the next step:

$$880 \phi_m(s) = \delta_m(s) + (1 - \delta_m(s))r_m(s) \quad (1)$$

- 881 • **Window** $W = (n, l)$: A block of l consecutive steps starting at step index n where all
882 states remain in \mathcal{U} . The ARM module replaces the baseline module m_{CoT} with a candidate
883 module m only on this window.
- 884 • **Visitation Weights** $w^\pi(W)$: The probability under baseline policy (π, m_{CoT}) that the
885 window W occurs. This measures the frequency with which the meta-policy starts the
886 module at a given window.

891 C.3 KEY SIMPLIFYING ASSUMPTIONS

892 The rest of our analysis relies on the following key assumptions.

893 **Assumption 1.** [Local competence lift in the scaffolding window] Within any given window $W =$
894 (n, l) , for all states s visited, the candidate model satisfies $\phi_m(s) \geq \phi_{m_{CoT}(s)} + \Delta_c$ for some lift
895 $\Delta_c \in [0, 1)$.

896 *Rationale:* This is the empirical premise that our scaffolded search objective (Section 3.2) is de-
897 signed to optimize for. Our Algorithm 1 directly measures and selects for modules that improve
898 local validity and recovery rates over the CoT within a constrained context at random locations.

899 **Assumption 2.** [Compatibility Loss] We define $\beta_l(W)$ a bound on the probability that replacing the
900 baseline m_{CoT} with ARM module m at a window W yields a context which is unusable for the rest
901 of the baseline reasoning trace. We refer to $(1 - \beta_l(W))$ as the *compatibility factor*. Furthermore,
902 we define $\bar{\beta}_l \sup_{W \in \mathcal{W}_{\text{valid}}} \beta_l(W)$ as the supremum of incompatibility probabilities across all valid
903 windows, representing the worst-case incompatibility bound.

904 *Rationale:* Swapping the baseline module m_{CoT} with m can introduce a “*context drift*” or “*semantic*
905 *drift*” which could amplify at deployment time when the ARM module m is used through the
906 entire trajectory. Our approach minimizes this drift by two means: 1) the few shot examples of the
907 progress, as well as the provided partial progress acts as a powerful inductive bias to constrain the
908 next step to states that preserve usefulness, for example by adopting the same notation, logical con-
909 tinuity, etc. 2) the reviewer agent which proposes mutations to the module (starting from baseline
910 CoT) is prompted to generate modules which solve one step at a time starting from the given partial
911 progress.

912 **Assumption 3.** [Module-invariant termination] We assume that the reward is terminal, and is pro-
913 vided under the condition that the extracted answer matches the right answer. Furthermore, the last
914 step in the MDP is reserved for this extraction, which is considered module-invariant, i.e. both CoT
915 or any other module can do this final syntax step (equally) perfectly.

918 C.4 THEORETICAL GROUNDING FOR THE SCAFFOLDED STEP-GENERATOR SEARCH

919 The scaffolded objective evaluates a candidate m by *splicing* it into a baseline rollout for a short
920 window $t \in \{i, \dots, i + \ell - 1\}$ while keeping m_{CoT} before and after:

$$921 \underbrace{U_{m_{CoT}}^* \circ (U_m^\ell) \circ U_{m_{CoT}}^i}_{\text{“baseline-candidate-baseline”}}.$$

922 This section formalizes the link between local module improvements and global performance gains.
923 Under our simplified MDP framework and assumptions 1, 2, we establish the following lemmas:
924

925 **Lemma 1.** [Window lift from local competence] The probability of remaining in \mathcal{U} after the window
926 increases by at least $l.C^{\ell-1}.\Delta_c$ for some constant $C \in (0, 1)$.
927

928 *Proof:* The l step window survival (being in \mathcal{U}) is at least $\phi_m(s)^l$. Under Assumption 1, the per-step
929 composite validity improves by at least Δ_c . Hence, the survival probability improves by at least
930 $(\phi_{CoT}(s) + \Delta_c)^l - \phi_{CoT}(s)^l$. Applying the mean value theorem on $f(x) = x^\ell$, we get

$$931 (x + \Delta_c)^\ell - x^\ell = \ell \xi^{\ell-1} \Delta_c \geq \ell \left(\min_{s \in \mathcal{U}_R} \phi_{CoT}(s) \right)^{\ell-1} \Delta_c$$

932 for some $\xi \in [x, x + \Delta]$, where \mathcal{U}_R represents states reachable by baseline CoT module. Further,
933 let $C := \min_{s \in \mathcal{U}_R} \phi_{CoT}(s)$ be a constant.
934

935 **Lemma 2.** [Accounting for Compatibility] The probability of a sample surviving the window W
936 and remain usable after the module swap is lower bounded by $(1 - \beta_l(W)).(l.C^{\ell-1}.\Delta_c)$ where C is
937 the constant from Lemma 1.

938 *Proof:* From Assumption 2, the usability probability is at least $(1 - \beta_l(W))$. Multiplying this by the
939 probability of surviving the window from Lemma 1 yields the result.
940

941 **Lemma 3.** [From window survival to finite-horizon success] Any increase in the probability of
942 staying within \mathcal{U} across a window W (while remaining usable) weakly increases the probability of
943 reaching \mathcal{G} within the horizon N .
944

945 *Proof:* Under assumption 3, the termination rule is module-invariant, and reaching the goal state
946 only depends on being in state $s \in \mathcal{U} \cup \mathcal{G}$ and $\text{solved}(s)=1$. Thus a higher probability of
947 preserving valid, in-progress states across a window cannot decrease (and generally increases) the
948 likelihood that subsequent steps will generate such a solvable state before the horizon is exhausted.
949 This follows from standard monotonicity arguments on absorbing Markov chains.
950

951 **Theorem 1** (Gain from Scaffolded Module Substitution in recursive meta policy). Let
952 $J(\pi_{Rec}, m) = \mathbb{E}[V_N^{\pi_{Rec}, m}(s_0)]$ denote the expected terminal reward (success probability) obtained
953 when recursively applying the step-generator module m under a fixed baseline meta-policy π_{Rec} for
954 horizon N . The improvement of the ARM module m^* over m_{CoT} is at least:
955

$$956 J(\pi_{Rec}, m^*) - J(\pi_{Rec}, m_{CoT}) \geq \sum_W w_\pi(W) \kappa(W) (1 - \beta_l(W)) l C^{\ell-1} \Delta_c. \quad (2)$$

957 where each term represents:

- 958 • $w_\pi(W)$: visitation probability of a window W under a baseline rollout;
- 959 • $\kappa(W)$: probability that a usable post-window state leads to terminal success within the
960 remaining horizon.
- 961 • $C \in (0, 1]$: a constant from Lemma 1 capturing compounding survival over steps.

962 In particular, if $\kappa(W) \geq \kappa_{min} \geq 0$ for all W , then

$$963 J(\pi_{Rec}, m^*) - J(\pi_{Rec}, m_{CoT}) \geq \kappa_{min} (1 - \beta_l(W)) l C^{\ell-1} \sum_W w_\pi(W) \quad (3)$$

964 *Proof:* From Lemma 2, the probability of surviving the window is lower bounded by $(1 - \beta_l)lC^{\ell-1}$.
965 Let $\kappa(W)$ represent the probability of success upon starting from a good state, post the window. By
966
967
968
969
970
971

Lemma 3, the increase in usable post-window mass translates to atleast a $\kappa(W)$ fraction improvement in terminal success within the remaining horizon. Thus the expected gain from the window is $\kappa(W)(1 - \beta_l(W))lC^{l-1}$. Taking expectation over window visitation probabilities yields the result:

$$\begin{aligned} J(\pi_{Rec}, m^*) - J(\pi_{Rec}, m_{CoT}) &= \sum_W w_\pi(W) \cdot \text{Gain}(W) \\ &= \sum_W w_\pi(W) \cdot \kappa(W) (1 - \beta_l(W)) l C^{l-1}. \\ &\geq \sum_W w_\pi(W) \cdot \kappa_{min} (1 - \beta_l(W)) l C^{l-1}. \end{aligned} \quad (4)$$

where $\kappa_{min} \geq 0$ is the lowest probability of success from a valid, usable intermediate reasoning trace.

C.5 THEORETICAL JUSTIFICATION FOR ZERO-SHOT POLICY TRANSFER

The learned meta policy π^* uses m_{CoT} as the *step generator* during the learning phase and is deployed zero-shot using the discovered ARM module m^* . Below, we justify why this transfer is effective.

Theorem 2 (Validity of Zero-shot step generator swap in Meta policy). Let $J(\pi^*, m^*) = \mathbb{E}_{s_0 \sim \mathcal{D}} [V_N^{(\pi^*, m^*)}(s_0)]$ denote the expected terminal reward obtained when applying the discovered meta-policy π^* (from Section 3.3), with the step-generator module m^* under horizon N . If $\Delta_c \geq \frac{\beta_l}{1-\beta_l}$, then the transfer is valid, i.e., $J(\pi^*, m^*) \geq J(\pi^*, m_{CoT})$

Let’s define per-step advantage of a module m over m_{CoT} with n more steps to go as the expected difference in value when taking one step with m and the rest with m_{CoT} :

$$A_n(s_n, m) \triangleq \mathbb{E}_{s' \sim P_m(\cdot|s)} [V_{n-1}^{(\pi^*, m_{CoT})}(s')] - \mathbb{E}_{s' \sim P_{m_{CoT}}(\cdot|s)} [V_{n-1}^{(\pi^*, m_{CoT})}(s')] \quad (5)$$

Now let’s consider the difference in expected value starting from a given state s_0 sampled from the data distribution \mathcal{D} . For simplicity, we drop π^* from notation as it is the common meta policy in both terms.

$$V_N^m(s_0) - V_N^{m_{CoT}}(s_0)$$

Rolling out for one step yields

$$\mathbb{E}_{s_1 \sim P_m(\cdot|s_0)} [V_{n-1}^m(s_1)] - V_n^{m_{CoT}}(s_0)$$

Adding an subtracting $\mathbb{E}_{s_1 \sim P_m(\cdot|s_0)} [V_{n-1}^{m_{CoT}}(s_1)]$ (i.e., sampling from m but continue with m_{CoT}) we get:

$$\begin{aligned} &\mathbb{E}_{s_1 \sim P_m(\cdot|s_0)} [V_{n-1}^{m_{CoT}}(s_1)] - \mathbb{E}_{s_1 \sim P_{m_{CoT}}(\cdot|s_0)} [V_{n-1}^{m_{CoT}}(s_1)] \\ &+ \mathbb{E}_{s_1 \sim P(\cdot|s_0)} [V_{n-1}^m(s_1) - V_{n-1}^{m_{CoT}}(s_1)] \end{aligned} \quad (6)$$

By Equation-5, this can be written as

$$A_n^{m_{CoT}}(s_0, m) + \mathbb{E}_{s_1 \sim P(\cdot|s_0)} [V_{n-1}^m(s_1) - V_{n-1}^{m_{CoT}}(s_1)]$$

This is a recursive equation in n since the second term is the difference in value between the module with $n - 1$ steps to horizon. Hence:

$$V_N^m(s_0) - V_N^{m_{CoT}}(s_0) = \sum_{n=0}^N [A_{N-n}^{m_{CoT}}(s_n, m)] \quad (7)$$

Thus we can conservatively guarantee module improvement, if each of the the advantage term is positive. Suppose that U represents the event that one step rollout using our discovered module m^* is usable (i.e. no errors, and usable context) in the next turn, then by law of total expectation the advantage term can be written as:

$$\begin{aligned} & \mathbb{P}(U | s, m^*) \cdot (\mathbb{E}[V_{N-n} | s, m^*, U] - \mathbb{E}[V_{N-n} | s, m_{\text{CoT}}, U]) \\ & + \mathbb{P}(\neg U | s, m^*) \cdot (\mathbb{E}[V_{N-n} | s, m^*, \neg U] - \mathbb{E}[V_{N-n} | s, m_{\text{CoT}}, \neg U]) \end{aligned} \quad (8)$$

By Assumption 1 and Assumption 2, the first term is at least $(1 - \bar{\beta}_l) \cdot \Delta_c$. The second term is lower bounded in the worst case by $\bar{\beta}_l \cdot (-1)$ since the probability of non-useful state is $\bar{\beta}_l$ and the difference in reward is at most -1 (when m^* fails where m_{CoT} succeeds). Thus a conservative lower bound is:

$$A_n^{m_{\text{CoT}}}(s, m) \geq (1 - \bar{\beta}_l) \cdot \Delta_c - \bar{\beta}_l \quad (9)$$

Hence advantage is non-negative when $\Delta_c \geq \frac{\bar{\beta}_l}{1 - \bar{\beta}_l}$, and we can guarantee that zero-shot transfer is effective when this condition is satisfied.

Corollary 1. The meta policy, paired with ARM, (π^*, m^*) yields the highest performance when Δ_c is large and β_l is small.

Remark: From Equation-8, the claim holds. Our scaffolded objective function is precisely designed to maximize Δ_c while reducing the probability of state drift leading to mismatched compatibility. In 3, we can empirically see that better performance is associated with improvement in per-step competence. And in 2, we see that the ARM module is compatible with m_{CoT} and can act as a superior drop-in replacement for CoT. Furthermore, the examples shown in Section D show that both CoT and ARM rollouts are empirically non-discriminable.

D STEP GENERATOR ROLL OUT EXAMPLES

The scaffolded objective function (Section 3.2) deliberately ensures that the discovered ARM module (m^*) is a strict drop-in **refinement** of the CoT operator rather than a qualitatively new operator. This ensures that compute scaling is performed at granular reasoning level, while learning interpretable, CoT like behaviors while reducing error propagation (Section 6.1& Section C). As seen in the following examples, the final rollout output of the ARM module closely resembles that of the CoT operator.

D.1 EXAMPLE-1

ARM Rollout

```
### Step 2
Express the sides in terms of the common ratio \ (2R\):
\[
a = 2R \sin 60^\circ = 2R \times \frac{\sqrt{3}}{2} = R \sqrt{3},
\]
\[
b = 2R \sin 45^\circ = 2R \times \frac{\sqrt{2}}{2} = R \sqrt{2},
\]
\[
c = 2R \sin 75^\circ.
\]
Recall that \ (\sin 75^\circ = \sin(45^\circ + 30^\circ)\ ):
\[
\sin 75^\circ = \sin 45^\circ \cos 30^\circ + \cos 45^\circ \sin 30^\circ
= \frac{\sqrt{2}}{2} \times \frac{\sqrt{3}}{2} + \frac{\sqrt{2}}{2} \times \frac{1}{2} = \frac{\sqrt{6}}{4} + \frac{\sqrt{2}}{4}
\]
\]
Hence,
\[
c = 2R \times \frac{\sqrt{6} + \sqrt{2}}{4} = R \times \frac{\sqrt{6} + \sqrt{2}}{2}.
\]
\]
```

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

CoT Rollout

```

### Step 2
Expressing the sides in terms of the circumradius  $(R)$ :
\[
a = 2R \sin A = 2R \sin 60^\circ = 2R \times \frac{\sqrt{3}}{2} = R \sqrt{3},
\]
\[
b = 2R \sin B = 2R \times \frac{\sqrt{2}}{2} = R \sqrt{2},
\]
\[
c = 2R \sin C = 2R \times \sin 75^\circ.
\]
Recall that:
\[
\sin 75^\circ = \sin (45^\circ + 30^\circ) = \sin 45^\circ \cos 30^\circ + \cos 45^\circ \sin 30^\circ = \frac{\sqrt{2}}{2} \times \frac{\sqrt{3}}{2} + \frac{\sqrt{2}}{2} \times \frac{1}{2} = \frac{\sqrt{6}}{4} + \frac{\sqrt{2}}{4} = \frac{\sqrt{6} + \sqrt{2}}{4}.
\]
Therefore:
\[
c = 2R \times \frac{\sqrt{6} + \sqrt{2}}{4} = R \times \frac{\sqrt{6} + \sqrt{2}}{2}.
\]

```

ARM Rollout

```

### Step 3
Now, determine the new temperature after doubling the kinetic energy:
\[
T_2 = 526.3 \text{ K}
\]
Convert this back to Celsius:
\[
T_2^\circ \text{ C} = 526.3 - 273.15 = 253.15^\circ \text{ C}
\]
This is a significant increase from the initial  $(-10^\circ \text{ C})$ , indicating the molecules now have enough energy corresponding to about  $(253^\circ \text{ C})$ .

```

CoT Rollout

```

### Step 3
Determine the new temperature after doubling the kinetic energy:
Since kinetic energy  $(\propto T)$ ,
\[
T_2 = 2 \times T_1 = 526.3 \text{ K}
\]
Convert this back to Celsius:
\[
T_2^\circ \text{ C} = T_2 \text{ (K)} - 273.15 = 526.3 - 273.15 \approx 253.15^\circ \text{ C}
\]
This temperature is well above the boiling point of water at 1 atm (100C). Therefore, the molecules' kinetic energy would be sufficient to overcome intermolecular forces and convert the water from the solid phase (ice) into the gas phase (vapor).

```

1134 D.3 EXAMPLE-3

1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

ARM Rollout

Step 1

To determine the correct substitution from equation 12.9 to equation 12.10, we need to understand the context of these equations and the variables involved. Equation 12.9 likely involves a momentum-like variable (P^+) or a gamma matrix that acts on spinors, and the goal is to simplify or rewrite this expression into a form involving gamma matrices or energy-momentum components.

Typically, in relativistic quantum mechanics or quantum field theory, the transformation from the light-cone momentum (P^+) to gamma matrices involves choosing a specific gamma matrix to project or decompose spinor components. The options provided suggest that the substitution involves gamma matrices and possibly the momentum operator or gamma matrices related to time.

Given that, the known structures are:

- (P^+) : light-cone momentum component, usually $(P^+ = P^0 + P^3)$.
- (γ^u) : a gamma matrix with a specific index (u) , possibly a spatial index.
- (γ^0) : the gamma matrix corresponding to time.
- $(\gamma^0 P)$: a product of (γ^0) with some momentum operator (P) .

Since the transformation is from equation 12.9 to 12.10, and the options involve gamma matrices and momenta, the typical approach in such derivations involves rewriting the projection operator or the momentum in terms of gamma matrices.

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

CoT

Step 1

To determine the correct substitution, first understand the context: equation 12.9 involves the momentum operator (P^μ) , and equation 12.10 is likely a transformed version involving gamma matrices. In relativistic quantum mechanics, transformations often involve gamma matrices (γ^μ) and the gamma matrices' role in formulating Dirac equations.

Given the options:

- A: (P^+)
- B: (γ^u)
- C: (γ^0)
- D: $(\gamma^0 P)$

and based on the typical transformations, the substitution from (P^μ) to a form involving gamma matrices usually involves (γ^0) because it relates to the time component of the gamma matrices and appears frequently in the Dirac formalism.

In particular, if the original equation involves (P^μ) , and the transformed equation involves a gamma matrix, then the key is recognizing the common substitution that converts momentum operators into gamma matrix expressions, especially the time component (P^0) .

Therefore, the step is to analyze which gamma matrix substitution correctly transforms the operator form in equation 12.9 into the one in 12.10.

E BEST ARM DISCOVERED: CRITICCHAINOFTHOUGHTV7

The following is the Python implementation of the best ARM discovered by our algorithm.

```

1242 1 import asyncio
1243 2
1244 3 class CriticChainOfThoughtV7:
1245 4     def __init__(self, llm):
1246 5         self.llm = llm
1247 6
1248 7     async def forward(self, problem, partial_progress):
1249 8         # 1. Generate four candidate next steps in parallel
1250 9         candidate_tasks = [
1251 10             self.llm.generate_step(problem, partial_progress)
1252 11             for _ in range(4)
1253 12         ]
1254 13         candidates = await asyncio.gather(*candidate_tasks)
1255 14
1256 15         # 2. Critique candidates in two groups of two, in parallel
1257 16         critique_tasks = []
1258 17         groups = [
1259 18             (0, 2, ("rating_1", "rating_2"), ("critique_1",
1260 19             ↪ "critique_2")),
1261 20             (2, 4, ("rating_3", "rating_4"), ("critique_3",
1262 21             ↪ "critique_4"))
1263 22         ]
1264 23         for start, end, rating_names, critique_names in groups:
1265 24             context = [
1266 25                 {
1267 26                     "name": "Problem",
1268 27                     "data": problem,
1269 28                     "description": "The problem to solve."
1270 29                 },
1271 30                 {
1272 31                     "name": "Partial Progress",
1273 32                     "data": partial_progress,
1274 33                     "description": "The partial solution so far."
1275 34                 },
1276 35                 {
1277 36                     "name": "Candidate Next Steps",
1278 37                     "data": "\n\n".join(
1279 38                         f"### Candidate Next Step
1280 39                         ↪ {i+1}\n{candidates[i]}"
1281 40                         for i in range(start, end)
1282 41                     ),
1283 42                     "description": "Two candidate next steps
1284 43                     ↪ formatted with markdown subheaders."
1285 44                 }
1286 45             ]
1287 46             instructions = (
1288 47                 "You are given a problem, the current partial
1289 48                 ↪ solution, and two candidate next reasoning steps.\n"
1290 49                 "For each candidate, provide:\n"
1291 50                 f"- {rating_names[0]} and {rating_names[1]}: a single
1292 51                 ↪ integer rating from 1 to 10 indicating its fit as the next
1293 52                 ↪ reasoning step (10 is best).\n"
1294 53                 f"- {critique_names[0]} and {critique_names[1]}: a
1295 54                 ↪ one-sentence critique highlighting each candidate's strengths
1296 55                 ↪ and weaknesses.\n"
1297 56                 f"Name the fields exactly {rating_names[0]},
1298 57                 ↪ {critique_names[0]}, {rating_names[1]}, {critique_names[1]}."
1299 58             )
1300 59             response_format = [
1301 60                 {

```

```

1296 51         "name": rating_names[0],
1297 52         "description": f"Integer rating (1-10) for
1298 ↪ Candidate Next Step {start+1}."
1299 53     },
1300 54     {
1301 55         "name": critique_names[0],
1302 56         "description": f"One-sentence critique of
1303 ↪ Candidate Next Step {start+1}."
1304 57     },
1305 58     {
1306 59         "name": rating_names[1],
1307 60         "description": f"Integer rating (1-10) for
1308 ↪ Candidate Next Step {start+2}."
1309 61     },
1309 62     {
1310 63         "name": critique_names[1],
1311 64         "description": f"One-sentence critique of
1312 ↪ Candidate Next Step {start+2}."
1313 65     }
1314 66 ]
1315 67 critique_tasks.append(
1316 68     self.llm.chat_completion(context, instructions,
1317 ↪ response_format)
1318 69 )
1319 70
1320 71 critiques = await asyncio.gather(*critique_tasks)
1321 72
1322 73 # 3. Parse ratings and identify the two highest-rated
1323 ↪ candidates
1324 74 ratings = [
1325 75     int(critiques[0]["rating_1"]),
1326 76     int(critiques[0]["rating_2"]),
1327 77     int(critiques[1]["rating_3"]),
1328 78     int(critiques[1]["rating_4"])
1329 79 ]
1330 80 sorted_indices = sorted(range(4), key=lambda i: ratings[i],
1331 ↪ reverse=True)
1332 81 top1_idx, top2_idx = sorted_indices[0], sorted_indices[1]
1333 82 top1_candidate = candidates[top1_idx]
1334 83 top2_candidate = candidates[top2_idx]
1335 84
1336 85 # 4. Final head-to-head comparison between the top two
1337 ↪ candidates
1338 86 context_final = [
1339 87     {
1340 88         "name": "Problem",
1341 89         "data": problem,
1342 90         "description": "The problem to solve."
1343 91     },
1344 92     {
1345 93         "name": "Partial Progress",
1346 94         "data": partial_progress,
1347 95         "description": "The partial solution so far."
1348 96     },
1349 97     {
1350 98         "name": "Candidate Next Steps",
1351 99         "data": (
1352 100             f"### Candidate A\n{top1_candidate}\n\n"
1353 101             f"### Candidate B\n{top2_candidate}"
1354 102         ),
1355 103         "description": "Two top candidate next steps
1356 ↪ formatted with markdown subheaders."
1357 104     }
1358 105 ]
1359 106 instructions_final = (

```

```

1350 107         "Compare Candidate A and Candidate B as the next
1351 108         ↪ reasoning step for the given problem and partial progress.\n"
1352 109         "Provide:\n"
1353 109         "- winner: choose either 'Candidate A' or 'Candidate B'
1354 110         ↪ indicating which step is better.\n"
1355 110         "- justification: one-sentence justification for your
1356 111         ↪ choice."
1357 111     )
1358 112     response_format_final = [
1359 113         {
1360 114             "name": "winner",
1361 115             "description": "Either 'Candidate A' or 'Candidate B'
1362 116         ↪ indicating the better next step."
1363 117         },
1364 118         {
1365 119             "name": "justification",
1366 120             "description": "One-sentence justification for the
1367 121         ↪ choice."
1368 121     }
1369 122     ]
1370 122     final_decision = await self.llm.chat_completion(
1371 123         context_final, instructions_final, response_format_final
1372 124     )
1373 125
1374 126     if final_decision["winner"].strip() == "Candidate A":
1375 127         selected_candidate = top1_candidate
1376 128         runnerup_candidate = top2_candidate
1377 129     else:
1378 130         selected_candidate = top2_candidate
1379 131         runnerup_candidate = top1_candidate
1380 132
1381 133     # 5. Post-selection adversarial critique with severity rating
1382 134     context_flaw = [
1383 135         {
1384 136             "name": "Problem",
1385 137             "data": problem,
1386 138             "description": "The problem to solve."
1387 139         },
1388 140         {
1389 141             "name": "Partial Progress",
1390 142             "data": partial_progress,
1391 143             "description": "The partial solution so far."
1392 144         },
1393 145         {
1394 146             "name": "Selected Candidate Next Step",
1395 147             "data": f"### Selected Candidate Next
1396 148         ↪ Step\n(selected_candidate)",
1397 149             "description": "The final chosen candidate next
1398 150         ↪ reasoning step formatted with a markdown subheader."
1399 150         }
1400 151     ]
1401 151     instructions_flaw = (
1402 152         "You are given a problem, the current partial solution,
1403 153         ↪ and a selected next reasoning step.\n"
1404 154         "Identify any major flaw or missing piece of reasoning in
1405 155         ↪ the selected step.\n"
1406 156         "Provide:\n"
1407 157         "- flaw: either the single word 'None' if there is no
1408 158         ↪ flaw, or a brief description of the flaw.\n"
1409 159         "- severity: a single integer rating from 1 to 10
1410 160         ↪ indicating how severe the flaw is (10 is critical)."
```

```

1404 161         "description": "Either the single word 'None' if
1405     ↪ there is no flaw, or a brief description of a major flaw in
1406     ↪ the selected step."
1407 162     },
1408 163     {
1409 164         "name": "severity",
1410 165         "description": "Integer rating (1-10) indicating
     ↪ severity of the flaw (10 is most severe)."

```

Listing 1: Code for CriticChainOfThoughtV7, performance: 38.0

F BEST META-POLICY DISCOVERED: VERIFIEDWEIGHTEDADAPTIVESELFCONSISTENTCHAINOFTHOUGHT

The following is the Python implementation of the best meta-policy discovered by our algorithm.

```

1438 1 import asyncio
1439 2 from agent.solution import Solution, Step
1440 3 from judge_utils import judge_equality
1441 4
1442 5 class VerifiedWeightedAdaptiveSelfConsistentChainOfThought:
1443 6     def __init__(self, llm, block):
1444 7         self.llm = llm
1445 8         self.block = block
1446 9
1447 10     async def forward(self, problem):
1448 11         # Helper: generate one chain up to 8 steps, then complete via
     ↪ LLM if needed
1449 12         async def generate_chain():
1450 13             solution = Solution()
1451 14             for _ in range(8):
1452 15                 next_step = await self.block.forward(problem,
     ↪ str(solution))
1453 16                 solution.add_step(Step(str(next_step)))
1454 17                 if solution.is_completed():
1455 18                     return solution
1456 19             completion = await self.llm.complete_solution(problem,
     ↪ str(solution))
1457 20             solution.add_step(Step(str(completion)))
1458 21             return solution

```

```

1458 22
1459 23     # Helper: confidence scoring (1-5)
1460 24     async def score_chain(chain):
1461 25         context = [
1462 26             {"name": "Problem", "data": problem, "description":
1463 27             ↪ "The original problem statement."},
1464 28             {"name": "Chain", "data": str(chain),
1465 29             ↪ "description": "Full chain-of-thought reasoning plus final
1466 30             ↪ answer."}
1467 31         ]
1468 32         instructions = (
1469 33             "You are evaluating the chain-of-thought solution for
1470 34             ↪ the given problem. "
1471 35             "On a scale from 1 (very uncertain) to 5 (very
1472 36             ↪ confident), rate your confidence "
1473 37             "that the final answer is correct. Output ONLY the
1474 38             ↪ integer confidence (1-5).")
1475 39         )
1476 40         response_format = [{"name": "Confidence", "description":
1477 41             ↪ "Integer from 1 to 5"}]
1478 42         resp = await self.llm.chat_completion(context,
1479 43             ↪ instructions, response_format)
1480 44         # parse safely
1481 45         try:
1482 46             conf = int(resp["Confidence"].strip())
1483 47         except Exception:
1484 48             conf = 1
1485 49         return max(1, min(conf, 5))
1486 50
1487 51     # Helper: verify logical consistency (Yes/No)
1488 52     async def verify_chain(chain):
1489 53         context = [
1490 54             {"name": "Problem", "data": problem, "description":
1491 55             ↪ "The original problem statement."},
1492 56             {"name": "Chain", "data": str(chain),
1493 57             ↪ "description": "Full chain-of-thought reasoning plus final
1494 58             ↪ answer."}
1495 59         ]
1496 60         instructions = (
1497 61             "Review the chain-of-thought reasoning for the given
1498 62             ↪ problem. "
1499 63             "Is the reasoning free of logical errors or
1500 64             ↪ contradictions? "
1501 65             "Output ONLY 'Yes' if it is fully logical, otherwise
1502 66             ↪ output 'No'."
1503 67         )
1504 68         response_format = [{"name": "Valid", "description": "Yes
1505 69             ↪ or No"}]
1506 70         resp = await self.llm.chat_completion(context,
1507 71             ↪ instructions, response_format)
1508 72         valid = resp.get("Valid",
1509 73             ↪ "").strip().lower().startswith("y")
1510 74         return valid
1511 75
1512 76     # Weighted vote helper
1513 77     def find_best_weighted(chains_list, conf_list):
1514 78         weight_sums = {}
1515 79         total = sum(conf_list)
1516 80         for chain, cf in zip(chains_list, conf_list):
1517 81             ans = chain.answer()
1518 82             weight_sums[ans] = weight_sums.get(ans, 0) + cf
1519 83         best_ans, best_w = None, -1
1520 84         for ans, w in weight_sums.items():
1521 85             if w > best_w:
1522 86                 best_ans, best_w = ans, w

```

```

1512 70         return best_ans, best_w, total
1513 71
1514 72         # 1) Generate initial 3 chains in parallel
1515 73         initial = [generate_chain() for _ in range(3)]
1516 74         chains = await asyncio.gather(*initial)
1517 75
1518 76         # 2) Score and verify each chain
1519 77         score_tasks = [score_chain(ch) for ch in chains]
1520 78         verify_tasks = [verify_chain(ch) for ch in chains]
1521 79         confidences = await asyncio.gather(*score_tasks)
1522 80         valids = await asyncio.gather(*verify_tasks)
1523 81
1524 82         max_chains = 7
1525 83
1526 84         # 3) Adaptive sampling with verification gating
1527 85         while True:
1528 86             # Determine which chains to consider: only verified if
1529 87             ↪ any, else all
1530 88             if any(valids):
1531 89                 considered_chains = [ch for ch, v in zip(chains,
1532 90                 ↪ valids) if v]
1533 91                 considered_confs = [cf for cf, v in
1534 92                 ↪ zip(confidences, valids) if v]
1535 93             else:
1536 94                 considered_chains = chains
1537 95                 considered_confs = confidences
1538 96
1539 97                 best_ans, best_weight, total_weight =
1540 98                 ↪ find_best_weighted(considered_chains, considered_confs)
1541 99                 # stop if weighted majority reached or chain cap
1542 100                 if best_weight > total_weight / 2 or len(chains) >=
1543 101                 ↪ max_chains:
1544 102                     break
1545 103
1546 104                 # else generate one more chain, score & verify, then loop
1547 105                 new_chain = await generate_chain()
1548 106                 chains.append(new_chain)
1549 107                 new_conf = await score_chain(new_chain)
1550 108                 confidences.append(new_conf)
1551 109                 new_valid = await verify_chain(new_chain)
1552 110                 valids.append(new_valid)
1553 111
1554 112         # 4) Select final chain: consensus & highest confidence among
1555 113         ↪ considered
1556 114         if any(valids):
1557 115             final_pool = [ (ch, cf) for ch, cf, v in zip(chains,
1558 116             ↪ confidences, valids) if v and judge_equality(ch.answer(),
1559 117             ↪ best_ans) ]
1560 118         else:
1561 119             final_pool = [ (ch, cf) for ch, cf in zip(chains,
1562 120             ↪ confidences) if judge_equality(ch.answer(), best_ans) ]
1563 121
1564 122         selected_chain = None
1565 123         top_conf = -1
1566 124         for ch, cf in final_pool:
1567 125             if cf > top_conf:
1568 126                 selected_chain, top_conf = ch, cf
1569 127
1570 128         # Fallback if nothing selected
1571 129         if selected_chain is None:
1572 130             selected_chain = chains[-1]
1573 131
1574 132         return selected_chain

```

Listing 2: Code for VerifiedWeightedAdaptiveSelfConsistentChainOfThought, performance: 38.0

G REPRODUCIBILITY STATEMENT

Upon publication, we commit to releasing the open-source code for our framework, including all discovered Agentic Reasoning Modules, meta-policies, and the specific prompts used for the Reviewer Agent. Our experiments were conducted using a mix of closed and open-source models. The MAS designer utilized OpenAI’s `o4-mini-high`. The reasoning modules were executed on GPT-4.1-nano, GPT-4o, and the open-source Llama-3.3-70B. All evaluation benchmarks, including MATH500, AIME, and HMMT, are publicly available.

G.1 ARM IMPLEMENTATION DETAILS

The 1000-sample subset of Open-R1-Mixture-of-Thoughts was created by taking the math and science splits of the original dataset, filtering to samples which the provided Deepseek-R1 reasoning trace had length between 8k to 10k tokens (to filter to samples of appropriate difficulty), and randomly sampling 1000 problems from the filtered problems.

We run both the step-generator module optimization and the meta-policy optimization for 20 iterations. Both optimizations are performed using GPT-4.1-nano as the MAS executor model.

Whenever sampling from the MAS executor model, we use a temperature of 1.0 with a `top_p` of 0.95.

G.2 BASELINE IMPLEMENTATION DETAILS

As in the ARM implementation, whenever sampling from the MAS executor model, we use a temperature of 1.0 with a `top_p` of 0.95.

- CoT: We use a simple CoT prompt that instructs the model to reason step-by-step and follow the final answer format.
- CoT-SC: We use $n = 12$ parallel reasoning traces.
- Self-Refine: We limit to a maximum of 5 self refining iterations.
- LLM-Debate: We use 4 LLM agents debating for a maximum of 3 rounds.
- ADAS: We use the provided codebase, following the recommended run configuration. For a fair comparison to other baselines, we make a one line addition to the optimizer prompt to disallow arbitrary Python code execution within the discovered MASes, since other baselines do not utilize code execution. For the 1000-sample optimization, we use GPT-4.1-nano as the MAS executor model during optimization, following ARM’s implementation.
- AFlow: We use the provided codebase, following the recommended run configuration. We allow the optimizer to utilize the Custom, AnswerGenerate, and ScEnsemble operators. For the 1000-sample optimization, we use GPT-4.1-nano as the MAS executor model during optimization, following ARM’s implementation.

H EXAMPLE OF REVIEWER FEEDBACK AND MUTATION

Here we show a concrete example of a reviewer surgical mutation, starting from an intermediate ARM module `CriticChainOfThoughtV2`, the reviewer feedback based on reflection on two rollouts, and corrective fixes to yield `CriticChainOfThoughtV3`.

ARM Mutation Example (`CriticChainOfThoughtV2`->`CriticChainOfThoughtV3`)

```

--- CriticChainOfThoughtV2.py
+++ CriticChainOfThoughtV3.py
import asyncio

class CriticChainOfThoughtV2:
    def __init__(self, llm):
        self.llm = llm

```

```

1620
1621
1622     async def forward(self, problem, partial_progress):
1623         # 1. Generate four candidate next steps in parallel
1624         candidate_tasks = [
1625             self.llm.generate_step(problem, partial_progress)
1626             for _ in range(4)
1627         ]
1628         candidates = await asyncio.gather(*candidate_tasks)
1629
1630         # 2. Critique candidates in two groups of two, in parallel
1631         ...
1632         # 3. Parse ratings and identify the two highest-rated
1633         ↪ candidates
1634         ...
1635         # 4. Final head-to-head comparison between the top two
1636         ↪ candidates
1637         ...
1638         # 5. Return the selected candidate
1639         - if final_decision["winner"].strip() == "Candidate A":
1640         -     return top1_candidate
1641         - else:
1642         -     return top2_candidate
1643         + if final_decision["winner"].strip() == "Candidate A":
1644         +     selected_candidate = top1_candidate
1645         +     runnerup_candidate = top2_candidate
1646         + else:
1647         +     selected_candidate = top2_candidate
1648         +     runnerup_candidate = top1_candidate
1649         +
1650         # 5. Post-selection adversarial critique of the chosen
1651         ↪ step
1652         context_flow = [
1653         +         {
1654         +             "name": "Problem",
1655         +             "data": problem,
1656         +             "description": "The problem to solve."
1657         +         },
1658         +         {
1659         +             "name": "Partial Progress",
1660         +             "data": partial_progress,
1661         +             "description": "The partial solution so far."
1662         +         },
1663         +         {
1664         +             "name": "Selected Candidate Next Step",
1665         +             "data": f"### Selected Candidate Next
1666         ↪ Step\n{selected_candidate}",
1667         +             "description": "The final chosen candidate next
1668         ↪ reasoning step formatted with a markdown subheader."
1669         +         }
1670         +     ]
1671         instructions_flow = (
1672         +         "You are given a problem, the current partial
1673         ↪ solution, and a selected next reasoning step.\n"
1674         +         "Identify any major flaw or missing piece of reasoning
1675         ↪ in the selected step.\n"
1676         +         "Provide:\n"
1677         +         "- flaw: either the single word 'None' if there is no
1678         ↪ flaw, or a brief description of the flaw."
1679         +     )
1680         response_format_flow = [
1681         +         {
1682         +             "name": "flaw",

```

```
1674
1675 +             "description": "Either the single word 'None' if
1676 ↪ there is no flaw, or a brief description of a major flaw in the
1677 ↪ selected step."
1678 +         }
1679 +     ]
1680 +     flaw_response = await self.llm.chat_completion(
1681 +         context_flaw, instructions_flaw, response_format_flaw
1682 +     )
1683 +     flaw = flaw_response["flaw"].strip()
1684 +     # 6. If a flaw is detected, fall back to the runner-up;
1685 ↪ otherwise keep the selected candidate
1686 +     if flaw.lower() != "none":
1687 +         return runnerup_candidate
1688 +     return selected_candidate
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
```

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

Reviewer Feedback on CriticChainOfThoughtV2

Performance Review of CriticChainOfThought CriticChainOfThoughtV2

Overall Performance Change

- CriticChainOfThought achieved an average score of 36.1, while CriticChainOfThoughtV2 reached 37.1 (an absolute gain of +1.0).
- This gain indicates that the added head-to-head comparison step in V2 yields a measurable improvement in next-step selection accuracy without disrupting baseline reasoning quality.

Evidence from the Execution Traces

1. Trace 1 (Divisibility Problem)

Both the previous block and V2 produced virtually identical Step 1 analyses, correctly identifying the need to enforce divisibility and positivity constraints. V2's head-to-head mechanism did not reduce the clarity or correctness of the initial reasoning step, demonstrating that the extra selection stage does not degrade generation quality.

2. Trace 2 (Complex Numbers, Roots of Unity)

Both blocks ultimately concluded that the set $\{S\}$ must be the $\{n\}$ th roots of unity for $\{n=3,4,5\}$, and both enclosed the final answer in `<final_answer>`. V2 preserved the previous blocks correct reasoning while adding no spurious steps, confirming that the new comparison stage retains core functionality.

3. Trace 3 (Consecutive Hits Combinatorics)

The previous blocks initial steps vacillated between interpreting 6.8 as 6, 7, or 8 shots across multiple retries. In contrast, V2's Step 1 consistently assumes 8 shots and clearly restates the pattern constraints before counting, suggesting that the head-to-head comparison helped V2 converge on a single, more accurate interpretation of the problem statement.

How the Evolution Drove the Improvement

- The original block generated diverse candidates and used two-way critiques to pick a top-rated step, but sometimes retained ambiguous or less precise interpretations.
- By introducing a third, final head-to-head comparison between the two frontrunners, V2 more reliably resolves borderline cases such as exactly how many total shots to assume in the combinatorics problem leading to slightly better overall accuracy.
- Crucially, this additional decision stage does not interfere with the correctness of initial reasoning steps or final answers, as seen in Traces 1 and 2, but does reduce interpretation noise in tricky examples like Trace 3.

Conclusion

CriticChainOfThoughtV2's addition of a head-to-head decision step raises next-step selection accuracy by resolving close calls between top candidates, yielding a +1.0 performance lift evident in more consistent problem interpretations while fully preserving the original blocks reasoning quality.

Reviewer Feedback on CriticChainOfThoughtV3

Performance Summary

The core change from CriticChainOfThoughtV2 to V3 was the addition of a post selection adversarial critique: after choosing between the top two candidate steps, V3 asks is there any major flaw? and, if the critic finds one, it falls back to the runner-up. In trace 1 (the quantum-mechanics problem) we see that

- V2 stopped at Step 3 without producing the required final answer (violating the if complete, immediately give <final_answer> rule).
- V3 both produced the correct reasoning step and appended the final answer `<final_answer>`.

However, over the full benchmark of challenging reasoning problems, V3's overall next-step selection accuracy dipped slightly (from 37.1% to 36.7%). This suggests that while the added adversarial critique can catch some subtle mistakes (and in the one provided trace corrected V2's omission of the final answer), it also sometimes overrules a correct choice leading to more false negatives than V2.

Key Observations

1. Completeness: V3 correctly obeyed the complete final answer rule in trace 1; V2 did not.
2. Average Accuracy: Despite the extra safeguard, V3's head-to-head + adversarial fallback pipeline results in a small net decrease in average next-step accuracy.
3. Critique Calibration: The adversarial critic occasionally identifies flaws where there are none, triggering unnecessary fallbacks.

Recommendations for V4

Loosen the fallback criterion or better calibrate the critics prompt so it only flags genuine errors. Consider logging how often the post-selection critique actually flips the choice, and analyze those cases to refine the critique instructions. Explore integrating a lighterweight self-check earlier (e.g. during rating) rather than as a hard post-selection veto.

I SEED CoT IMPLEMENTATION

The following is the Python implementation of the seed CoT step-generator that we use as the seed module and the corresponding prompt.

```

1 class ChainOfThought:
2     def __init__(self, llm):
3         self.llm = llm
4
5     async def forward(self, problem, partial_progress):
6         generated_step = await self.llm.generate_step(problem,
7             ↪ partial_progress)
8         return generated_step

```

Listing 3: Code for CoT Seed

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

```
CoT Prompt (Step-Generator)

## SYSTEM TURN
You will be tasked to complete a partial solution to a problem.

Continue the partial solution in a step by step manner, denoting the
  first step with "### Step {n}" and so on.

Each step should significantly progress the partial solution, without
  repeating any work done in previous steps.

Once the solution is completed, your response should conclude by
  giving the final answer enclosed within <final_answer> </
  final_answer>.

Do not include units in the final answer. For multiple choice
  questions, the final answer should only be the letter of the
  correct answer choice.

If the partial solution is already complete, immediately provide the
  final answer instead of writing another step.

## USER TURN
Problem:

{problem}

Complete this partial solution to the problem:

{partial_progress}
```