

# Learn Like Humans: Use Meta-cognitive Reflection for Efficient Self-Improvement

Anonymous ACL submission

## Abstract

While Large Language Models (LLMs) enable complex autonomous behavior, current agents remain constrained by static, human-designed prompts that limit adaptability. Existing self-improving frameworks attempt to bridge this gap but typically rely on inefficient, multi-turn recursive loops that incur high computational costs. To address this, we propose **Metacognitive Agent Reflective Self-improvement (MARS)**, a framework that achieves efficient self-evolution within a single recurrence cycle. Inspired by educational psychology, MARS mimics human learning by integrating principle-based reflection (abstracting normative rules to avoid errors) and procedural reflection (deriving step-by-step strategies for success). By synthesizing these insights into optimized instructions, MARS allows agents to systematically refine their reasoning logic without continuous online feedback. Extensive experiments on six benchmarks demonstrate that MARS outperforms state-of-the-art self-evolving systems while significantly reducing computational overhead. Code are available at <https://anonymous.4open.science/r/MARS-9F16>

## 1 Introduction

Large language models (LLMs) have enabled autonomous agents capable of complex reasoning, planning, and tool use (Brown et al., 2020; Wei et al., 2022; Yao et al., 2023). However, current agents rely on fixed, human-designed components—manually crafted prompts, predefined workflows, and static configurations—limiting their adaptability to strategies within human intuition (Hu et al., 2025a; Wang et al., 2024). While machine learning history shows hand-designed solutions are consistently replaced by learned ones (Elsken et al., 2019; Zoph and Le, 2017), agent development remains largely manual. The theoretical basis for self-improving AI dates back to Schmidhuber’s

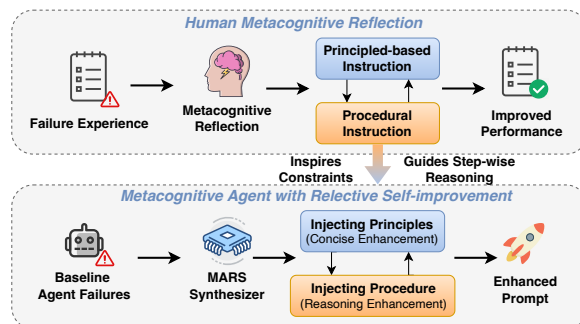


Figure 1: The Cognitive Inspiration behind MARS. This framework parallels human reflection with the MARS (Metacognitive Agent with Reflective Self-improvement) agent, converting baseline agent failures into principled-based and procedural instructions to synthesize enhanced prompts.

Gödel machines (Schmidhuber, 2007), which formalized self-referential systems that rewrite their own code. Although formal proof requirements make the original framework impractical (Steunebrink and Schmidhuber, 2011), it has inspired modern self-improving systems that rely on empirical validation instead.

However, current self-improvement frameworks for LLM agents tend to be constrained by multi-turn recursiveness, which results in inefficient learning and adaptation, as well as excessive computational resource usage. Humans, by contrast, are able to resolve previous errors and adapt to new solutions more efficiently through structured learning approaches. Research in education science has identified two complementary paradigms for guiding learners (Hiebert and Lefevre, 1986; Anderson, 1983). The first is *principle-based learning*, which focuses on helping learners avoid mistakes by establishing conceptual categories of what is correct versus incorrect, and understanding the underlying rules that govern a domain (Hiebert and Lefevre, 1986; Rittle-Johnson et al., 2001). The second is *procedural learning*, which emphasizes

067 using prior experience and step-by-step reasoning  
068 to increase the likelihood of successful outcomes  
069 (Anderson, 1983; Kolb, 1984). Rather than learn-  
070 ing in isolation, humans benefit most when they  
071 integrate both approaches through systematic re-  
072 flection and summarization of their experiences.  
073 Studies in metacognition have shown that struc-  
074 tured reflection—where learners explicitly analyze  
075 what worked, what failed, and why—significantly  
076 improves learning efficiency and knowledge trans-  
077 fer (Flavell, 1979; Kaplan et al., 2013; Stanton  
078 et al., 2021). Furthermore, research on productive  
079 failure demonstrates that learning from one’s own  
080 errors, when properly guided, leads to deeper con-  
081 ceptual understanding than direct instruction alone  
082 (Kapur, 2014, 2010).

083 In this work, we propose **MARS** (Metacognitive  
084 Agent with Reflective Self-improvement), a frame-  
085 work that enables multi-agent systems to achieve  
086 efficient self-improvement within a single recur-  
087 rence cycle by integrating both principle-based and  
088 procedural learning approaches. Inspired by human  
089 metacognitive learning, MARS allows agents to  
090 systematically reflect on their experiences, extract-  
091 ing general principles that help avoid past mistakes  
092 while simultaneously deriving procedural knowl-  
093 edge that replicates successful strategies. Unlike  
094 existing self-evolving agent frameworks that rely  
095 on multi-turn recursive improvement, which often  
096 leads to inefficient learning and excessive computa-  
097 tional costs, MARS consolidates the learning pro-  
098 cess through structured summarization, enabling  
099 agents to maximize adaptation efficiency in each  
100 improvement cycle.

101 Our main contributions are as follows:

- 102 • We propose MARS, a self-improvement  
103 framework for multi-agent systems that inte-  
104 grates principle-based and procedural learning  
105 inspired by human meta-cognitive theory.
- 106 • We introduce a triple-pathway reflection  
107 mechanism that extracts: (1) normative princi-  
108 ples for error avoidance, (2) procedural strate-  
109 gies for success replication, and (3) a unified  
110 synthesis of both pathways.
- 111 • We design a structured summarization mod-  
112 ule that consolidates learning within a single  
113 cycle, reducing computational overhead from  
114 multi-turn recursive improvement.
- 115 • We conduct extensive experiments on chal-  
116 lenging knowledge and reasoning bench-

marks, showing MARS outperforms exist- 117  
ing self-evolving frameworks while requiring 118  
fewer iterations. 119

## 2 Related Work 120

Recent research has transitioned from static 121  
prompting to *self-evolving agents*—systems capa- 122  
ble of analyzing their own performance, learning 123  
from errors, and modifying their behavior to im- 124  
prove over time. Drawing from meta-learning prin- 125  
ciples (Finn et al., 2017; Hospedales et al., 2022), 126  
these approaches can be broadly categorized into 127  
two paradigms: verbal reflection and structural self- 128  
modification. 129

The first category utilizes verbal reflection to 130  
facilitate learning from failure. *Reflexion* (Shinn 131  
et al., 2023) introduced this paradigm by enabling 132  
agents to generate natural language critiques of 133  
their mistakes, storing them in episodic memory 134  
to guide future reasoning. *RISE* (Qu et al., 2024) 135  
extends this by training models to iteratively detect 136  
and correct errors across multiple turns, demon- 137  
strating that self-correction capabilities can be in- 138  
ternalized through fine-tuning. While effective, 139  
these methods primarily rely on inference-time re- 140  
cursion or memory retrieval rather than permanent 141  
parameter or prompt optimization. The second 142  
category focuses on automated architecture and 143  
code evolution. Systems like *ADAS* (Hu et al., 144  
2025a) employ meta-agents to iteratively gener- 145  
ate and evaluate new agent designs in code, while 146  
*AgentSquare* (Shang et al., 2025) adopts a modular 147  
approach to evolve components for planning, rea- 148  
soning, and tool use. Similarly, *Agent-Pro* (Zhang 149  
et al., 2024) optimizes policies through reflection 150  
on historical trajectories. Taking this further, fully 151  
self-referential approaches allow agents to modify 152  
their own underlying source code. The *Gödel Agent* 153  
(Yin et al., 2025) enables agents to rewrite their 154  
logic guided by high-level objectives, a concept ex- 155  
tended by the *Darwin Gödel Machine* (Zhang et al., 156  
2025) and *SICA* (Robeyns et al., 2025), which in- 157  
tegrate evolutionary search to explore diverse self- 158  
improvement paths. 159

Despite these advancements, current frameworks 160  
face significant efficiency bottlenecks. Reflection- 161  
based methods often depend on computationally 162  
expensive multi-turn recursive loops, while code- 163  
modifying agents require complex validation envi- 164  
ronments. Unlike these approaches, our framework 165  
draws inspiration from human metacognitive the- 166

ory to achieve efficient self-improvement within a *single recurrence cycle*.

### 3 Methodology

To achieve efficient self-improvement without the computational overhead of recursive loops, we propose MARS, a three-phase framework designed to systematically transform sporadic model failures into targeted, actionable prompt enhancements. Rather than treating errors as isolated incidents, our approach aggregates failures to identify systematic weaknesses, synthesizes remediation strategies, and integrates them into a self-improving loop. Figure 2 illustrates the complete pipeline.

The framework operates as follows. In the *Evaluation* phase, an analyzer model  $\mathcal{M}_\phi$  examines each failed question and produces a structured analysis  $\mathcal{A}_i$  capturing both question characteristics (type  $\tau_i$ , topics  $\mathcal{T}_i$ ) and failure attributes (error type  $\epsilon_i$ , root cause  $\rho_i$ , specific mistake  $\mu_i$ ). The *Failure Allocation* phase then applies a grouping function  $\kappa$  to partition analyses into groups  $\mathcal{G} = \{G_j\}$  based on shared type-topic keys, aggregating diagnostic attributes into group-level error profiles  $\Psi_j$ . Finally, the *Enhancement Generation* phase synthesizes targeted enhancements  $(E_j^{(c)}, E_j^{(r)})$  for each group and combines them with the base prompt  $P$  via weighted aggregation to produce enhanced prompts  $P^{(c)}$  and  $P^{(r)}$ . We detail each phase below.

#### 3.1 Evaluation

The first phase of our enhancement pipeline performs fine-grained diagnosis of each incorrectly answered question. Rather than treating failures as a homogeneous set, we analyze each instance independently to understand the precise reasoning breakdown that led to the incorrect response.

Formally, let  $\mathcal{Q} = \{q_i\}_{i=1}^n$  denote a set of failed questions from the benchmark evaluation. Each instance  $q_i$  comprises question text, options, ground-truth answer  $a_i^*$ , the model’s predicted answer  $\hat{a}_i$ , and generated reasoning trace. We employ a specialized analyzer model,  $\mathcal{M}_\phi$ , to dissect these components. For every  $q_i$ , the analyzer produces a structured analysis

$$\mathcal{A}_i = (\tau_i, \mathcal{T}_i, \epsilon_i, \rho_i, \mu_i) \quad (1)$$

, which encapsulates two distinct categories of attributes:

The first category characterizes question itself. It assigns a “question type”  $\tau_i \in \mathcal{Y}$  (where  $\mathcal{Y} =$

{factual, conceptual, calculation, application}) and identifies a set of “topics”  $\mathcal{T}_i \subseteq \mathcal{D}$  derived from domain vocabulary  $\mathcal{D}$ . As detailed in Table 1, the combination of question type  $\tau_i$  and topic  $\mathcal{T}_i$  serves as composite key for grouping failures in the subsequent Allocation phase. The second category characterizes specific error mechanism. This includes an “error type”  $\epsilon_i \in \mathcal{E}$ , a natural language “root cause”  $\rho_i$  explaining the fundamental reasoning deficit, and a “specific mistake”  $\mu_i$  pinpointing the exact step where logic diverged. The error taxonomy  $\mathcal{E}$ , presented in Table 2, defines six categories ranging from conceptual misunderstandings to calculation errors.

To ensure consistent classification, we enforce a strict exclusivity rule: each failure is assigned to exactly one category in  $\mathcal{E}$ . In cases where multiple failure modes co-occur (e.g., a calculation error stemming from a conceptual misunderstanding), the analyzer assigns the category corresponding to the *earliest point of divergence* in the reasoning chain. The final output of this phase is the collection of structured analyses  $\mathbb{A} = \{\mathcal{A}_i\}_{i=1}^n$ .

#### 3.2 Failure Allocation

This aggregation transforms sparse, per-instance observations into dense, group-level patterns. By clustering failures with shared characteristics, the allocation phase enables the subsequent system to generate high-level guidance that addresses classes of errors simultaneously, ensuring that the enhanced prompts are both targeted and scalable.

The second phase organizes individual error analyses into semantically coherent groups to enable pattern discovery. While the previous phase examined each failure in isolation, this phase identifies structural similarities across failures that may share common remediation strategies. We define a composite grouping function

$$\kappa : \mathbb{A} \rightarrow \mathcal{Y} \times 2^{\mathcal{D}}, \quad \kappa(\mathcal{A}_i) = (\tau_i, \mathcal{T}_i) \quad (2)$$

that maps each analysis to its “type-topic key” via  $\kappa(\mathcal{A}_i) = (\tau_i, \mathcal{T}_i)$ . This two-dimensional grouping captures the intuition that errors on calculation questions about thermodynamics likely stem from different causes than errors on conceptual questions about molecular biology, and thus require distinct enhancement strategies.

Given the analysis set  $\mathbb{A}$  from the previous phase, we construct a partition

$$\Psi_j = (\mathcal{E}_j, \mathcal{R}_j, \mathcal{F}_j) \quad (3)$$

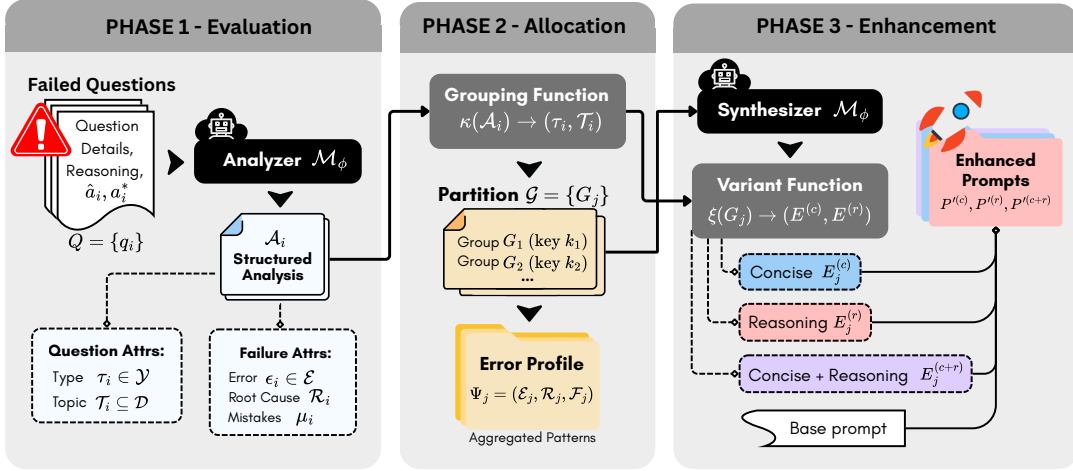


Figure 2: Overview of the proposed framework: (1) diagnose failed questions into structured analyses  $\mathcal{A}_i$ , (2) group by type-topic keys and aggregate error profiles  $\Psi_j$ , and (3) synthesize enhancements via weighted aggregation to produce  $P^{(c)}$  and  $P^{(r)}$ .

Type	Description	Example
Factual	Recall of specific facts or definitions	“What is the atomic number of carbon?”
Conceptual	Understanding of principles or theories	“Why does entropy increase in isolated systems?”
Calculation	Quantitative problem requiring computation	“Calculate the final velocity given...”
Application	Applying knowledge to novel scenarios	“Which treatment would be most effective for...”
Analysis	Breaking down complex information	“What can be inferred from this experimental data?”
Comparison	Evaluating similarities or differences	“Which compound has higher boiling point?”

Table 1: Question type categories.

, where each group  $G_j = \{\mathcal{A}_i \in \mathbb{A} : \kappa(\mathcal{A}_i) = k_j\}$  contains all analyses sharing the same type-topic key  $k_j$ . Within each group, we aggregate the diagnostic attributes to form a collective “error profile”  $\Psi_j = (\mathcal{E}_j, \mathcal{R}_j, \mathcal{F}_j)$ , comprising the set of observed error types  $\mathcal{E}_j = \{\epsilon_i : \mathcal{A}_i \in G_j\}$ , recurring root causes  $\mathcal{R}_j = \{\rho_i : \mathcal{A}_i \in G_j\}$ , and common difficulty factors  $\mathcal{F}_j$ .

This aggregation transforms sparse, per-instance observations into dense, group-level patterns. By clustering failures with shared characteristics, the allocation phase enables the subsequent system to generate high-level guidance that addresses classes of errors simultaneously, ensuring that the enhanced prompts are both targeted and scalable.

### 3.3 Enhancement Generation

The final phase synthesizes the group-level error patterns into prompt enhancements that guide the model toward correct reasoning. For each group  $G_j \in \mathcal{G}$ , we first perform pattern analysis to extract actionable guidance, then integrate this guidance into the original prompt template.

Given a group  $G_j$  with its aggregated error pro-

file  $\Psi_j$ , we query the analyzer model  $\mathcal{M}_\phi$  to synthesize targeted remediation strategies. The analyzer examines the common error types  $\mathcal{E}_j$ , shared root causes  $\mathcal{R}_j$ , and recurring mistakes within the group, then produces structured guidance including typical pitfalls to avoid, verification steps to perform, and domain-specific reasoning strategies. This group-level synthesis captures patterns that may not be apparent from any single failure but emerge clearly when examining multiple related errors.

We define an enhancement generation function

$$\xi : \mathcal{G} \rightarrow \mathcal{S}^{(c)} \times \mathcal{S}^{(r)}, \quad \xi(G_j) = (E_j^{(c)}, E_j^{(r)}) \quad (4)$$

that produces two distinct enhancement variants to accommodate different reasoning scenarios. The “concise” enhancement  $E_j^{(c)} \in \mathcal{S}^{(c)}$  provides brief warnings and key points, suitable for quick reference during inference. The “reasoning” enhancement  $E_j^{(r)} \in \mathcal{S}^{(r)}$  supplies minimal hints designed to trigger self-correction without overconstraining the reasoning process. Formally, for each group  $G_j$ , we produce an enhancement pair  $\xi(G_j) = (E_j^{(c)}, E_j^{(r)})$ .

Category	Description	Typical Manifestation
Conceptual Misunderstanding	Fundamental confusion about domain principles	Misapplying laws; confusing related concepts
Calculation Error	Computational or mathematical mistakes	Arithmetic errors; unit conversion mistakes
Misreading	Misinterpretation of question or choices	Overlooking negation; misidentifying the question
Incomplete Analysis	Premature termination of reasoning	Stopping after partial solution
Wrong Elimination	Incorrect rejection of candidate answers	Eliminating correct answer based on flawed logic
Knowledge Gap	Absence of requisite domain knowledge	Missing facts; unfamiliar terminology

Table 2: Error categories and descriptions.

Benchmark	Focus	Categories		Details	
		Primary Domains	Sub-categories	Key Features	Size
DROP (Dua et al., 2019)	Discrete Reasoning	Numerical Operations	Add/Sub, Min/Max, Count, Select, Compare	NFL & History passages	96K
MGSM (Shi et al., 2023)	Multilingual Math	Grade-school Math	Arithmetic, Word Problems, Algebra	11 languages (BN, SW, TE incl.)	250
MMLU (Hendrycks et al., 2021)	General Knowledge	STEM Humanities Social Sciences	Physics, Chemistry, CS, Math, Biology History, Philosophy, Law, Ethics Economics, Psychology, Politics	57 subjects; Elem.–Prof.	15.9K
GPQA (Rein et al., 2024)	Graduate Science	Biology Physics Chemistry	Molecular Bio (8%), Genetics General (10%), Electromagnetism Organic (36%), General	PhD-level; Google-proof	448
HLE (Phan et al., 2025)	Expert Academic	Mathematics (41%) Sciences (27%) Tech & Humanities (32%)	Algebra, Analysis, Combinatorics Physics, Chemistry, Biology CS/AI, Engineering, Social Sci.	100+ areas; 14% multimodal	2.5K
Omni-MATH (Gao et al., 2024)	Olympiad Math	Algebra & Number Theory Geometry Analysis & Discrete	Linear, Abstract, Primes, Modular Euclidean, Analytic, Projective Calculus, Combinatorics, Graph Theory	33+ sub-domains; 10 levels	4.4K

Table 3: Question categories across six LLM evaluation benchmarks used for category-based hybrid enhancement.

The final enhanced prompt  $P'$  is constructed by appending the relevant enhancements to the base prompt  $P$ . We define an aggregation operator  $\oplus$  that combines enhancements weighted by group cardinality  $|G_j|$ , prioritizing guidance derived from larger groups where more failures share the same type-topic characteristics. The resulting enhanced prompts are given by

$$P^{(c)} = P \oplus \bigoplus_{j=1}^m w_j E_j^{(c)}, \quad P^{(r)} = P \oplus \bigoplus_{j=1}^m w_j E_j^{(r)} \quad (5)$$

, where  $w_j \propto |G_j|$ , each embedding the collective remediation knowledge extracted from the failure analysis pipeline.

Beyond applying enhancements uniformly, we introduce a **Hybrid** strategy that dynamically selects the optimal enhancement type per question category. We apply four strategies: Concise, Reasoning, Concise+Reasoning, and Hybrid. The first three serve as ablation baselines. For the Hybrid strategy, each dataset is partitioned into train, validation, and test splits (8:1:1). The training set generates enhancements via Phases 1-3. The validation set determines which enhancement type performs best for each category  $c$ :

$$E_c^* = \arg \max_{E \in \{E^{(c)}, E^{(r)}, E^{(c+r)}\}} \text{Acc}(E, \mathcal{V}_c) \quad (6)$$

where  $\mathcal{V}_c$  denotes validation questions in category  $c$  and  $E^{(c+r)}$  combines both enhancement types. The selected  $E_c^*$  is applied to matching test questions.

## 4 Experiments

**Datasets.** We evaluate on six benchmarks spanning reasoning capacity and knowledge coverage. For reasoning capacity, we use DROP (Dua et al., 2019), a reading comprehension benchmark requiring discrete reasoning operations such as addition, counting, and sorting; MGSM (Shi et al., 2023), the Multilingual Grade School Math benchmark containing 250 problems; and OMNI-math (Gao et al., 2024), an Olympiad-level benchmark with 4,428 competition problems across 33 sub-domains. For knowledge coverage, we use MMLU (Hendrycks et al., 2021), which spans 57 subjects including STEM, humanities, and social sciences; GPQA (Rein et al., 2024), a graduate-level “Google-proof” Q&A benchmark with expert-written questions in biology, physics, and chemistry; and HLE (Phan et al., 2025) (Humanity’s Last Exam), a frontier benchmark with 2,500 expert-level questions designed to test the limits of current AI systems.

For DROP, MGSM, MMLU, and GPQA, we use gpt-3.5-turbo for comparison with baselines. For the more challenging benchmarks, Omni-MATH and Humanity’s Last Exam, we use gpt-4o due to its stronger reasoning capabilities and more recent training data cutoff.

Method	Enhancement	Reasoning		Knowledge	
		DROP	MGSM	MMLU	GPQA
MetaAgentSearch (Hu et al., 2025b)	-	79.4 <sup>†</sup>	53.4 <sup>†</sup>	69.6 <sup>†</sup>	34.6 <sup>†</sup>
Gödel Agent (Yin et al., 2025)	-	80.9 <sup>†</sup>	64.2 <sup>†</sup>	70.9 <sup>†</sup>	34.9 <sup>†</sup>
Zero-shot (Brown et al., 2020)	n/a	62.0	35.0	64.0	11.8
	Concise	63.5 <sub>+1.5</sub>	37.6 <sub>+2.6</sub>	60.7 <sub>-3.3</sub>	11.8 <sub>+0.0</sub>
	Reasoning	65.2 <sub>+3.2</sub>	37.0 <sub>+2.0</sub>	64.0 <sub>+0.0</sub>	12.7 <sub>+0.9</sub>
	Concise+Reasoning	63.8 <sub>+1.8</sub>	35.2 <sub>+0.2</sub>	64.2 <sub>+0.2</sub>	12.7 <sub>+0.9</sub>
	Hybrid	<b>68.4</b> <sub>+6.4</sub>	<b>39.4</b> <sub>+4.4</sub>	<b>65.1</b> <sub>+1.1</sub>	<b>20.0</b> <sub>+8.2</sub>
Zero-shot-CoT (Kojima et al., 2022)	n/a	74.5	52.9	65.8	16.4
	Concise	77.2 <sub>+2.7</sub>	54.4 <sub>+1.5</sub>	66.8 <sub>+1.0</sub>	19.1 <sub>+2.7</sub>
	Reasoning	78.8 <sub>+4.3</sub>	54.1 <sub>+1.2</sub>	69.0 <sub>+3.2</sub>	19.1 <sub>+2.7</sub>
	Concise+Reasoning	78.1 <sub>+3.6</sub>	54.6 <sub>+1.7</sub>	69.0 <sub>+3.2</sub>	17.3 <sub>+0.9</sub>
	Hybrid	<b>81.6</b> <sub>+7.1</sub>	<b>56.4</b> <sub>+3.5</sub>	<b>70.5</b> <sub>+4.7</sub>	<b>22.2</b> <sub>+5.8</sub>
Self-Refine (Madaan et al., 2023)	n/a	77.8	57.5	48.8	<b>36.4</b>
	Concise	80.5 <sub>+2.7</sub>	60.5 <sub>+3.0</sub>	60.5 <sub>+11.7</sub>	<b>38.2</b> <sub>+1.8</sub>
	Reasoning	<b>82.1</b> <sub>+4.3</sub>	58.4 <sub>+0.9</sub>	59.0 <sub>+10.2</sub>	<b>40.9</b> <sub>+4.5</sub>
	Concise+Reasoning	<b>81.2</b> <sub>+3.4</sub>	58.7 <sub>+1.2</sub>	63.9 <sub>+15.1</sub>	32.7 <sub>-3.7</sub>
	Hybrid	<b>84.3</b> <sub>+6.5</sub>	<b>61.3</b> <sub>+3.8</sub>	<b>64.6</b> <sub>+15.8</sub>	<b>49.1</b> <sub>+12.7</sub>
Self-Consistency (Wang et al., 2023)	n/a	79.5	63.5	61.8	18.2
	Concise	<b>83.8</b> <sub>+4.3</sub>	<b>73.7</b> <sub>+10.2</sub>	69.3 <sub>+7.5</sub>	26.4 <sub>+8.2</sub>
	Reasoning	<b>84.5</b> <sub>+5.0</sub>	<b>73.2</b> <sub>+9.7</sub>	<b>71.3</b> <sub>+9.5</sub>	24.6 <sub>+6.4</sub>
	Concise+Reasoning	<b>83.2</b> <sub>+3.7</sub>	<b>73.7</b> <sub>+10.2</sub>	<b>71.7</b> <sub>+9.9</sub>	21.8 <sub>+3.6</sub>
	Hybrid	<b>86.2</b> <sub>+6.7</sub>	<b>74.3</b> <sub>+10.8</sub>	<b>72.5</b> <sub>+10.7</sub>	<b>33.6</b> <sub>+15.4</sub>

Table 4: Results on DROP, MGSM, MMLU, and GPQA benchmarks. DROP uses F1 score; others use accuracy (%). Subscripts indicate improvement over n/a baseline. **Bold** indicates results exceeding Gödel Agent. <sup>†</sup> are results obtained from Yin et al., 2025.

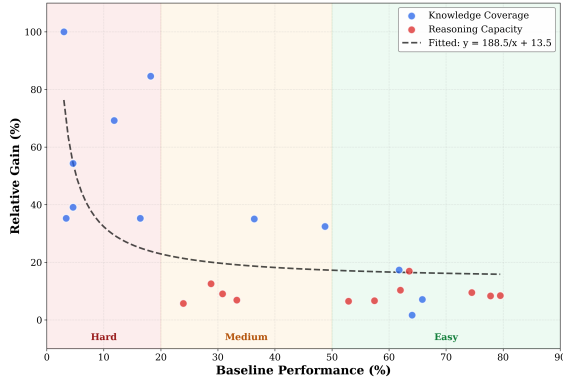
**Implementation.** We evaluate four base prompting methods: Zero-shot (Brown et al., 2020), Zero-shot-CoT (Kojima et al., 2022) which appends “Let’s think step by step” to elicit reasoning, Self-Refine (Madaan et al., 2023) which iteratively critiques and improves responses, and Self-Consistency (Wang et al., 2023) which samples multiple reasoning paths and selects answers via majority voting.

We apply four enhancement strategies to each prompting method: Concise (principle-based do’s and don’ts for avoiding common errors), Reasoning (explicit step-by-step instructions to follow correct rationale), Concise+Reasoning (combining both), and Hybrid (dynamically selecting strategies based on question category). Concise, Reasoning, and Concise+Reasoning serve as ablation baselines to isolate the contribution of each enhancement type. For the Hybrid strategy, we leverage the question categories outlined in Table 3: DROP (5 discrete reasoning types), MGSM (3 math types  $\times$  11 languages), MMLU (57 subjects across 4 domains), GPQA (3 scientific domains), Humanity’s Last Exam (8 broad categories), and Omni-MATH (33+ mathematical sub-domains). Each dataset is split into train:val:test = 8:1:1. The validation set is

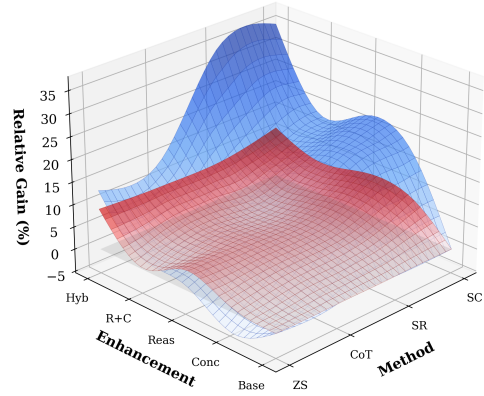
used to discover the optimal enhancement strategy for each question category, which is then applied to the corresponding categories in the test set. We compare against MetaAgentSearch and Gödel Agent (Yin et al., 2025), in which Gödel Agent represents a state-of-the-art meta-learning optimized agent system. All experiments use temperature  $T = 0$  and maximum token length of 3,000. For Self-Consistency, we sample  $n = 5$  responses with temperature  $T = 0.7$ . We report F1 score for DROP and accuracy for all other benchmarks.

## 5 Results and Analysis

**Main Results** Table 4 presents results across four benchmarks: reasoning capacity (DROP and MGSM) and knowledge coverage (MMLU and GPQA). We compare our enhancement strategies against MetaAgentSearch (Hu et al., 2025b) and Gödel Agent (Yin et al., 2025). MetaAgentSearch and Gödel Agent employ instruction-free self-improvement through multiple recursive iterations, automatically discovering agent architectures without human-crafted guidance. In contrast, methods like Self-Refine (Madaan et al., 2023) rely on



(a) Baseline performance vs. relative gain from hybrid enhancement. Background shading indicates difficulty tiers.



(b) 3D surface of relative performance gain across prompting methods (X), enhancement types (Y), and accuracy (Z).

Figure 3: Performance gain analysis. (a) Inverse relationship between task difficulty and enhancement effectiveness. (b) Performance landscape comparison between Knowledge Coverage (blue) and Reasoning Capacity (red).

Method	Enhancement	Reasoning	Knowledge
		OMNI-math	HLE
Zero-shot	n/a	23.93	3.40
	Concise	23.44 <sub>-0.49</sub>	3.20 <sub>0.20</sub>
	Reasoning	24.56 <sub>+0.63</sub>	3.40 <sub>+0.00</sub>
	R+C ★	22.43 <sub>-1.50</sub>	4.20 <sub>+0.80</sub>
	Hybrid	25.30 <sub>+1.37</sub>	4.60 <sub>+1.20</sub>
Zero-shot-CoT	n/a	30.81	4.60
	Concise	31.04 <sub>+0.23</sub>	5.10 <sub>+0.50</sub>
	Reasoning	31.83 <sub>+1.02</sub>	5.40 <sub>+0.80</sub>
	R+C ★	31.72 <sub>+0.91</sub>	5.20 <sub>+0.60</sub>
	Hybrid	33.60 <sub>+2.79</sub>	6.40 <sub>+1.80</sub>
Self-Refine	n/a	28.78	4.60
	Concise	28.67 <sub>-0.11</sub>	6.40 <sub>+1.80</sub>
	Reasoning	30.59 <sub>+1.81</sub>	5.50 <sub>+0.90</sub>
	R+C ★	29.12 <sub>+0.34</sub>	6.00 <sub>+1.40</sub>
	Hybrid	32.40 <sub>+3.62</sub>	<b>7.10</b> <sub>+2.50</sub>
Self-Consistency	n/a	33.30	3.00
	Concise	33.60 <sub>+0.30</sub>	3.40 <sub>+0.40</sub>
	Reasoning	34.60 <sub>+1.30</sub>	5.20 <sub>+2.20</sub>
	R+C ★	32.50 <sub>-0.80</sub>	4.80 <sub>+1.80</sub>
	Hybrid	<b>35.60</b> <sub>+2.30</sub>	6.00 <sub>+3.00</sub>

Table 5: Results on OMNI-math and HLE. Accuracy (%) reported. Subscripts indicate improvement over baseline. **Bold** indicates best result per dataset. ★s are short for reasoning plus concise enhancements

human-crafted prompts encoding explicit refinement strategies without groundtruth. Notably, even without our enhancements, Self-Refine (36.4%) already outperforms both Gödel Agent (34.9%) and MetaAgentSearch (34.6%) on GPQA. This demonstrates that well-designed human-crafted prompting can surpass instruction-free recursive optimization. With hybrid enhancement, Self-Consistency surpasses Gödel Agent on three benchmarks (DROP: 86.2 vs. 80.9, MGSM: 74.3 vs. 64.2, MMLU: 72.5 vs. 70.9), while Self-Refine with hy-

brid reaches 49.1% on GPQA. These results suggest that MARS offers a cost-effective alternative to complex recursive agent systems.

Zero-shot and Zero-shot-CoT show consistent improvements, with hybrid yielding +6.4 and +7.1 on DROP respectively. Self-Refine achieves dramatic gains on knowledge benchmarks: +15.8 on MMLU and +12.7 on GPQA. Self-Consistency benefits substantially across all benchmarks (+10.8 on MGSM, +10.7 on MMLU), indicating synergy between multiple reasoning paths and enhanced prompts. Hybrid enhancement consistently yields the largest improvements, validating category-aware selection. Reasoning generally outperforms Concise on reasoning-intensive tasks. Combining both (Concise+Reasoning) does not always yield additive benefits—on GPQA with Self-Refine, it underperforms individual enhancements (32.7% vs. 40.9% for Reasoning), suggesting interference when naively combining strategies.

**Generalization to Challenging Benchmarks** To investigate whether MARS generalizes to more challenging evaluation settings, we evaluate on two additional benchmarks: Omni-MATH and Humanity’s Last Exam (HLE). Table 5 presents these results.

Omni-MATH tests advanced mathematical reasoning with baseline accuracies below 35%. Despite this difficulty, our enhancements provide consistent improvements. Reasoning enhancement proves particularly effective, yielding gains across all prompting methods (+0.63 for Zero-shot, +1.02 for Zero-shot-CoT, +1.81 for Self-Refine, +1.30 for Self-Consistency). Hybrid enhancement achieves

the best results, with Self-Consistency reaching 35.60% (+2.30). Notably, Concise+Reasoning often underperforms individual enhancements (e.g., 32.50% vs. 34.60% for Self-Consistency), consistent with the interference pattern observed in main results.

HLE represents an extremely challenging benchmark with baseline accuracies below 5%. Even in this difficult regime where models operate near floor-level performance, MARS demonstrates effectiveness. Self-Refine with hybrid enhancement achieves 7.10%, a 54.3% relative improvement over baseline (4.60%). Zero-shot-CoT with hybrid reaches 6.40% (+1.80), and Self-Consistency with Reasoning achieves 5.20% (+2.20). These gains indicate that our category-aware enhancements extract meaningful improvements even on problems designed to challenge state-of-the-art systems. Finally, using a different model for enhancement generation did not result in significant changes, as presented in Appendix D. Qwen2.5-72B-Instruct-Turbo produced comparable enhancement patterns across both knowledge coverage and reasoning capacity benchmark that confirms MARS enhancement’s effectiveness generalizes across enhancement generators rather than being dependent on a specific model.

**Performance Gain Analysis** We analyze the relationship between baseline performance and relative gain from prompt enhancement using scatter plots and 3D surface visualizations across knowledge coverage (HLE, GPQA, MMLU) and reasoning capacity (OMNI-math, MGSM, DROP) benchmarks (Figure 3).

A significant inverse correlation exists between baseline performance and relative gain (Spearman  $\rho = -0.654$ ,  $p < 0.001$ ), with the fitted model gain =  $188.54/\text{baseline} + 13.48$  ( $R^2 = 0.443$ ). Critically, this relationship is category-dependent: knowledge coverage datasets exhibit strong correlation ( $\rho = -0.795$ ,  $p = 0.002$ ), while reasoning capacity datasets show no significant relationship ( $\rho = 0.264$ ,  $p = 0.433$ ). This divergence suggests fundamentally different enhancement mechanisms—knowledge tasks benefit disproportionately at low baselines, whereas reasoning gains remain uniform across difficulty levels.

The 3D surface analysis reveals a method-enhancement interaction specific to reasoning tasks. For reasoning datasets, self-consistency combined with one-turn MARS enhancement produces signif-

icantly amplified gains compared to other prompting methods (7.31% vs. 2.66%; Mann-Whitney  $U$ ,  $p = 0.050$ ). This amplification pattern is *not observed* in knowledge coverage datasets, where high variance ( $\sigma = 21.28\%$ ) obscures potential interaction effects and gains are primarily explained by the baseline-gain relationship.

Items	Statistic	$p$ -value
Overall baseline-gain correlation	$\rho = -0.654$	$< 0.001^{***}$
Knowledge baseline-gain correlation	$\rho = -0.795$	$= 0.002^{**}$
Reasoning baseline-gain correlation	$\rho = 0.264$	$= 0.433$
SC amplification (Reasoning)	$U$	$= 0.050^*$
Categories differ in gain	$U$	$= 0.003^{**}$

Table 6: Statistical significance summary for gain analysis.

These findings suggest category-aware enhancement strategies: for knowledge tasks, prioritize low-baseline scenarios where gains are maximized; for reasoning tasks, leverage self-consistency methods which uniquely amplify MARS enhancement effects.

## 6 Conclusion

We presented MARS, a metacognitive framework that integrates principle-based reflection (learning what to avoid) with procedural reflection (learning how to succeed) for efficient self-improvement in LLM agents. Existing instruction-free self-improvement methods rely on multi-turn recursive optimization that is both computationally expensive and often underperforming. MARS overcomes both shortcomings by consolidating learning into a single recurrence cycle through structured summarization, while generating targeted, category-aware enhancements from systematic failure analysis. Experiments across six benchmarks demonstrate that MARS consistently outperforms state-of-the-art self-evolving systems with significantly reduced computational overhead, suggesting that human-inspired learning paradigms offer a practical alternative to resource-intensive recursive self-improvement.

## Limitations

Two main limitations warrant consideration. First, the predefined error taxonomy may not generalize to all task types. Our six error categories and type-topic groupings are designed for structured benchmarks with clear correct answers. For open-ended or creative tasks where errors are less cate-

550	gorical and more subjective, these classifications	Timothy Hospedales, Antreas Antoniou, Paul Micaelli,	599
551	may prove insufficient. Extending MARS to such	and Ajo Storkey. 2022. <a href="#">Meta-learning in neural net-</a>	600
552	domains would require developing more flexible	<a href="#">works: A survey</a> . <i>IEEE Transactions on Pattern Anal-</i>	601
553	categorization schemes.	<i>ysis and Machine Intelligence</i> , 44(9):5149–5169.	602
554	Second, single-cycle learning trades depth for ef-	Shengran Hu, Cong Lu, and Jeff Clune. 2025a. <a href="#">Auto-</a>	603
555	ficiency. While our approach significantly reduces	<a href="#">mated design of agentic systems</a> . In <i>International</i>	604
556	computational cost compared to recursive methods,	<i>Conference on Learning Representations</i> .	605
557	it inherently limits the improvement achievable in	Shengran Hu, Cong Lu, and Jeff Clune. 2025b. <a href="#">Au-</a>	606
558	one pass. For applications prioritizing maximum	<a href="#">tomated design of agentic systems</a> . <i>Preprint</i> ,	607
559	performance over efficiency, we recommend apply-	arXiv:2408.08435.	608
560	ing MARS iteratively—using enhanced prompts	Matthew Kaplan, Naomi Silver, Danielle LaVaque-	609
561	to generate new failure sets, then deriving further	Manty, and Deborah Meizlish. 2013. <i>Using Reflec-</i>	610
562	refinements from residual errors.	<i>tion and Metacognition to Improve Student Learning:</i>	611
		<i>Across the Disciplines, Across the Academy</i> . Stylus	612
		Publishing, Sterling, VA.	613
563	<b>References</b>	Manu Kapur. 2010. Productive failure in mathematical	614
564	John R. Anderson. 1983. <i>The Architecture of Cognition</i> .	problem solving. <i>Instructional Science</i> , 38(6):523–	615
565	Harvard University Press, Cambridge, MA.	550.	616
566	Tom Brown, Benjamin Mann, Nick Ryder, Melanie	Manu Kapur. 2014. Productive failure in learning math.	617
567	Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind	<i>Cognitive Science</i> , 38(5):1008–1022.	618
568	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-	619
569	Askeell, and 1 others. 2020. Language models are few-	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-	620
570	shot learners. In <i>Advances in Neural Information</i>	guage models are zero-shot reasoners. In <i>Advances in</i>	621
571	<i>Processing Systems</i> , volume 33, pages 1877–1901.	<i>Neural Information Processing Systems</i> , volume 35,	622
572	Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel	pages 22199–22213.	623
573	Stanovsky, Sameer Singh, and Matt Gardner. 2019.	David A. Kolb. 1984. <i>Experiential Learning: Expe-</i>	624
574	<a href="#">DROP: A reading comprehension benchmark requir-</a>	<i>rience as the Source of Learning and Development</i> .	625
575	<a href="#">ing discrete reasoning over paragraphs</a> . In <i>Proceed-</i>	Prentice-Hall, Englewood Cliffs, NJ.	626
576	<i>ings of NAACL-HLT</i> .	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	627
577	Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter.	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	628
578	2019. Neural architecture search: A survey. <i>Journal</i>	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	629
579	<i>of Machine Learning Research</i> , 20(55):1–21.	Shashank Gupta, Bodhisattwa Prasad Majumder,	630
580	Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017.	Katherine Hermann, Sean Welleck, Amir Yazdan-	631
581	Model-agnostic meta-learning for fast adaptation of	bakhsh, and Peter Clark. 2023. <a href="#">Self-refine: Iterative</a>	632
582	deep networks. <i>International Conference on Ma-</i>	<a href="#">refinement with self-feedback</a> . In <i>Advances in Neu-</i>	633
583	<i>chine Learning</i> , pages 1126–1135.	<i>ral Information Processing Systems</i> , volume 36.	634
584	John H. Flavell. 1979. Metacognition and cognitive	Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li,	635
585	monitoring: A new area of cognitive-developmental	Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang,	636
586	inquiry. <i>American Psychologist</i> , 34(10):906–911.	Mohamed Shaaban, John Ling, Sean Shi, Michael	637
587	Bofei Gao and 1 others. 2024. <a href="#">Omni-math: A univer-</a>	Choi, Anish Agrawal, Arnav Chopra, Adam Khoja,	638
588	<a href="#">sal olympiad level mathematic benchmark for large</a>	Ryan Kim, Richard Ren, Jason Hausenloy, Oliver	639
589	<a href="#">language models</a> . <i>arXiv preprint arXiv:2410.07985</i> .	Zhang, Mantas Mazeika, and 317 others. 2025. <a href="#">Hu-</a>	640
590	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	<a href="#">manity’s last exam</a> . <i>Preprint</i> , arXiv:2501.14249.	641
591	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	Yuxiao Qu, Tianjun Zhang, Naman Garg, and Aviral	642
592	2021. <a href="#">Measuring massive multitask language under-</a>	Kumar. 2024. <a href="#">Recursive introspection: Teaching</a>	643
593	<a href="#">standing</a> . <i>ICLR</i> .	<a href="#">language model agents how to self-improve</a> . In <i>Ad-</i>	644
594	James Hiebert and Patricia Lefevre. 1986. Conceptual	<i>vances in Neural Information Processing Systems</i> ,	645
595	and procedural knowledge in mathematics: An in-	volume 37.	646
596	troductory analysis. In <i>Conceptual and Procedural</i>	Qwen Team. 2024. <a href="#">Qwen2.5 technical report</a> . <i>arXiv</i>	647
597	<i>Knowledge: The Case of Mathematics</i> , pages 1–27.	<i>preprint arXiv:2412.15115</i> .	648
598	Lawrence Erlbaum Associates, Hillsdale, NJ.	David Rein, Betty Li Hou, Asa Cooper Stickland,	649
		Jackson Petty, Richard Yuanzhe Pang, Julien Di-	650
		rani, Julian Michael, and Samuel R. Bowman. 2024.	651
		<a href="#">GPQA: A graduate-level google-proof q&amp;a bench-</a>	652
		<a href="#">mark</a> . <i>COLM</i> .	653

654	Bethany Rittle-Johnson, Robert S. Siegler, and Martha Wagner Alibali. 2001. Developing conceptual understanding and procedural skill in mathematics: An iterative process. <i>Journal of Educational Psychology</i> , 93(2):346–362.	707
655		708
656		709
657		710
658		711
659	Maxime Robeyns, Martin Szummer, and Laurence Aitchison. 2025. <a href="#">A self-improving coding agent</a> . <i>arXiv preprint arXiv:2504.15228</i> . ICLR 2025 Workshop on Scaling Self-Improving Foundation Models.	712
660		713
661		714
662		715
663	Jürgen Schmidhuber. 2007. <a href="#">Gödel machines: Fully self-referential optimal universal self-improvers</a> . In <i>Artificial General Intelligence</i> , pages 199–226. Springer, Berlin, Heidelberg.	716
664		717
665		718
666		719
667	Yu Shang, Yu Li, Keyu Zhao, Likai Ma, Jiahe Liu, Fengli Xu, and Yong Li. 2025. <a href="#">Agentsquare: Automatic llm agent search in modular design space</a> . In <i>International Conference on Learning Representations</i> .	720
668		721
669		722
670		723
671		724
672	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2023. <a href="#">Language models are multilingual chain-of-thought reasoners</a> . In <i>ICLR</i> .	725
673		726
674		727
675		728
676		729
677	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. <a href="#">Reflection: Language agents with verbal reinforcement learning</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 36.	730
678		731
679		732
680		733
681		734
682	Julie Dangremond Stanton, Amanda J. Sebesta, and John Dunlosky. 2021. Fostering metacognition to support student learning and performance. <i>CBE—Life Sciences Education</i> , 20(2):fe3.	735
683		736
684		737
685		738
686	Bas R. Steunebrink and Jürgen Schmidhuber. 2011. A family of gödel machine implementations. In <i>International Conference on Artificial General Intelligence</i> , pages 275–280. Springer.	739
687		740
688		741
689		742
690	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji-Rong Wen. 2024. <a href="#">A survey on large language model based autonomous agents</a> . <i>Frontiers of Computer Science</i> , 18(6).	743
691		744
692		745
693		746
694		747
695		748
696	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. <a href="#">Self-consistency improves chain of thought reasoning in language models</a> . In <i>International Conference on Learning Representations</i> .	749
697		750
698		751
699		752
700		753
701		754
702	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. <a href="#">Chain-of-thought prompting elicits reasoning in large language models</a> . <i>Advances in Neural Information Processing Systems</i> , 35.	755
703		756
704		757
705		758
706		759
	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. <a href="#">React: Synergizing reasoning and acting in language models</a> . <i>International Conference on Learning Representations</i> .	760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

---

**Algorithm 1** MARS Enhancement Pipeline

---

**Require:** Failed questions  $\mathcal{Q} = \{q_1, \dots, q_n\}$ , ground truths  $\{a_i^*\}$ , predictions  $\{\hat{a}_i\}$ , base prompt  $P$ , analyzer  $\mathcal{M}_\phi$ , error taxonomy  $\mathcal{E}$

**Ensure:** Enhanced prompts  $(P^{(c)}, P^{(s)}, P^{(r)})$

```
1:                                     ▷ Phase 1: Evaluation
2:  $\mathbb{A} \leftarrow \emptyset$ 
3: for  $q_i \in \mathcal{Q}$  do
4:    $\mathcal{A}_i \leftarrow \text{DIAGNOSE}(\mathcal{M}_\phi, q_i, a_i^*, \hat{a}_i)$ 
5:    $\mathbb{A} \leftarrow \mathbb{A} \cup \{\mathcal{A}_i\}$ 

6:                                     ▷ Phase 2: Failure Allocation
7:  $\mathcal{G} \leftarrow \text{CLUSTER}(\mathbb{A})$ 

8:                                     ▷ Phase 3: Enhancement Generation
9:  $\mathbb{E} \leftarrow \emptyset$ 
10: for  $G_j \in \mathcal{G}$  do
11:    $(E_j^{(c)}, E_j^{(s)}, E_j^{(r)}) \leftarrow \text{SYNTHESIZE}(\mathcal{M}_\phi, G_j)$ 
12:    $\mathbb{E} \leftarrow \mathbb{E} \cup \{(E_j^{(c)}, E_j^{(s)}, E_j^{(r)})\}$ 

13:                                     ▷ Aggregate into final prompts
14:  $(P^{(c)}, P^{(s)}, P^{(r)}) \leftarrow \text{AGGREGATE}(P, \mathbb{E}, \mathcal{G})$ 

15: return  $(P^{(c)}, P^{(s)}, P^{(r)})$ 
```

---

745 purpose during inference. Table 8 summarizes their  
746 characteristics.

747 The **concise** variant  $E^{(c)}$  provides brief warnings  
748 derived from common mistakes, suitable for scenarios  
749 where inference cost is a concern. The **specific**  
750 variant  $E^{(s)}$  includes detailed verification steps and  
751 explicit reasoning strategies, appropriate when ac-  
752 curacy is prioritized over efficiency. The **reasoning**  
753 variant  $E^{(r)}$  offers minimal guidance designed to  
754 trigger self-correction without over-constraining  
755 the model’s reasoning process.

## 756 C Prompts

757 This appendix provides the complete prompt tem-  
758 plates used in our experiments. All prompts  
759 use a structured XML-style output format with  
760 <reasoning> and <answer> tags to facilitate  
761 consistent response parsing. The placeholder  
762 {question} is replaced with the specific problem  
763 instance at inference time.

### 764 C.1 Zero-Shot Prompting

765 The zero-shot baseline provides the question di-  
766 rectly without any demonstrations or reasoning in-  
767 structions.

### 768 C.2 Zero-Shot Chain-of-Thought

769 Zero-shot chain-of-thought prompting (Kojima  
770 et al., 2022) elicits step-by-step reasoning without  
771 providing exemplars.

### 772 C.3 Few-Shot Chain-of-Thought

773 Few-shot chain-of-thought prompting (Wei et al.,  
774 2022) provides exemplars demonstrating the rea-  
775 soning process.

### 776 C.4 Self-Consistency

777 Self-consistency (Wang et al., 2023) samples mul-  
778 tiple reasoning paths and aggregates answers via  
779 majority voting. We evaluate configuration with 10  
780 samples.

### 781 C.5 Self-Refine

782 Self-refine (Madaan et al., 2023) enables iterative  
783 improvement of responses through self-feedback.

### 784 C.6 Enhancement Analyzer

785 The Enhancement Analyzer is the first LLM-  
786 powered agent in our zero-shot enhancement  
787 pipeline. It performs individual failure analysis by  
788 examining each incorrectly answered GPQA ques-  
789 tion to determine the precise cause of error. Given  
790 a failed question along with the model’s predicted  
791 answer and reasoning, the analyzer classifies the  
792 question type (factual, conceptual, calculation, ap-  
793 plication, analysis, or comparison), identifies spe-  
794 cific scientific topics, and diagnoses the error type  
795 (e.g., conceptual misunderstanding, calculation er-  
796 ror, misreading, incomplete analysis, wrong elimi-  
797 nation, or knowledge gap). The agent produces a  
798 structured JSON output containing the root cause  
799 explanation, the specific reasoning step that failed,  
800 required domain knowledge, and factors contribut-  
801 ing to the question’s difficulty. This granular analy-  
802 sis enables downstream pattern recognition across  
803 multiple failures.

### 804 C.7 Enhancement Synthesizer

805 The Enhancement Synthesizer is the second LLM-  
806 powered agent that operates on grouped failures  
807 sharing common characteristics. After the En-  
808 hancement Analyzer processes individual ques-  
809 tions, failures are clustered by question type and  
810 topic. The Synthesizer then analyzes each cluster to  
811 identify recurring error patterns, synthesize shared  
812 root causes, and generate targeted enhancement  
813 strategies. For each type-topic group, it produces  
814 common mistake patterns, critical warnings, verifi-  
815 cation steps, topic-specific guidance, and a concise  
816 prompt addition designed to prevent similar errors.  
817 This synthesized knowledge is then used to con-  
818 struct three variants of enhanced prompts (concise,  
819 specific, and reasoning-focused), each tailored to

Variant	Purpose	Content	Length
Concise ( $E^{(c)}$ )	Quick reference	Warnings and key points	Short
Specific ( $E^{(s)}$ )	Thorough guidance	Mistakes, verification steps, approach	Detailed
Reasoning ( $E^{(r)}$ )	Self-correction	Minimal hints for discovery	Minimal

Table 8: Comparison of three enhancement variants.

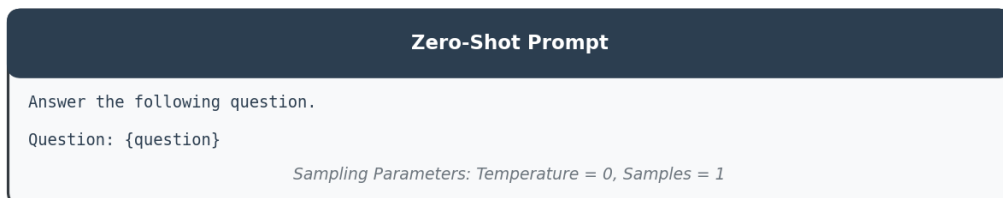


Figure 4: Zero-shot prompt template.

address the identified weaknesses while maintaining the base prompting strategy’s structure.

### C.8 Model Configuration

Table 9 summarizes the model configuration used across all experiments.

Parameter	Value
Default Model	GPT-4o
Evaluation Model	o3-mini
Max Tokens	2000
Timeout	30 seconds

Table 9: Model configuration settings.

Component	Concise	Reasoning	Specific
Explicit warnings	Yes	No	Yes
Action sequences	Yes	No	Yes
Process guidance	Brief	Detailed	Detailed
Verification steps	No	Yes	Yes
Approach description	No	Yes	Yes
Relative length	Short	Medium	Long

Table 10: Structural comparison of MARS enhancement types.

## D Enhancement Generation with Open-source Model

To evaluate the generalizability of MARS enhancements across different enhancement generators, we replicate experiments using Qwen2.5-72B-Instruct-Turbo (Qwen Team, 2024) for enhancement generation on one knowledge coverage dataset (GPQA) and one reasoning capacity dataset (OMNI-math). Results are presented in Tables 11 and 12.

**Analysis.** Qwen-generated enhancements demonstrate consistent improvement patterns across both datasets. The hybrid enhancement achieves the highest gains, with self-refine + hybrid reaching 48.20% on GPQA (a 32.6% relative improvement over baseline). Individual enhancements show modest gains: concise enhancement improves zero-shot by 7.7% relative, while reasoning+concise provides the largest single-enhancement gain for self-consistency (24.9% relative). These results confirm that MARS enhancement effectiveness generalizes across enhancement generators, though optimal enhancement selection remains task-dependent.

## E One-Turn MARS Enhancement Samples

MARS generates three types of one-turn prompt enhancements from the same error analysis, each with different structures and verbosity levels. We demonstrate each type using the Algebra\_Equations category under zero-shot prompting.

### E.1 Concise Enhancement

The concise enhancement (Figure 11) provides compact, action-oriented guidance organized by type-topic groups. Each group includes warning indicators ([!]) highlighting common failure patterns with failure counts, followed by action arrows (->) specifying recommended problem-solving procedures. This format minimizes prompt length while preserving critical guidance.

```

Zero-Shot Chain-of-Thought Prompt

Solve this problem step by step.
Question: {question}

Provide your response in the following format:

<reasoning>
[Show your step-by-step work and explain your thinking process here]
</reasoning>

<answer>
[Provide only the final answer here]
</answer>

Sampling Parameters: Temperature = 0, Samples = 1

```

Figure 5: Zero-shot chain-of-thought prompt template.

```

Few-Shot Chain-of-Thought Prompt

Here are some examples of how to solve similar problems:

Example 1:
Question: If  $x + 3 = 7$ , what is  $x$ ?
<reasoning>
To find  $x$ , I need to isolate it on one side.
Subtracting 3 from both sides:  $x + 3 - 3 = 7 - 3$ 
Simplifying:  $x = 4$ 
</reasoning>
<answer>
4
</answer>

Now solve this problem:
Question: {question}

Provide your response in the following format:

<reasoning>
[Show your step-by-step work and explain your thinking process here]
</reasoning>

<answer>
[Provide only the final answer here]
</answer>

Sampling Parameters: Temperature = 0, Samples = 1

```

Figure 6: Few-shot chain-of-thought prompt template with one demonstration example.

## E.2 Reasoning Enhancement

The reasoning enhancement (Figure 12) emphasizes problem-solving strategies over explicit warnings. Each type-topic group receives a bulleted consideration (\*) describing the recommended reasoning approach. This format focuses on *how* to think about problems rather than *what* to avoid, making it effective for tasks requiring flexible reasoning.

## E.3 Specific Enhancement (Reasoning + Concise)

The specific enhancement (Figure 13) combines both approaches into a comprehensive three-part structure for each type-topic group: (1) **Common Mistakes** (x) explicitly enumerate failure patterns;

(2) **Verification Steps** (+) provide concrete validation actions; (3) **Approach** describes the recommended problem-solving methodology. This format maximizes guidance completeness at the cost of increased prompt length.

## E.4 Enhancement Comparison

Table 10 summarizes the structural differences between enhancement types.

The hybrid enhancement strategy evaluates all three types on a validation set and selects the optimal enhancement for each question, combining their complementary strengths.

865  
866  
867  
868  
869  
870  
871  
872  
873  
  
874  
875  
  
876  
877  
878  
879

880  
881  
882  
883  
884  
  
885  
  
886  
887  
  
888  
889  
890  
891

```

Self-Consistency Prompt (10 Samples)

Solve this problem step by step.
Question: {question}

Provide your response in the following format:

<reasoning>
[Show your step-by-step work and explain your thinking process here]
</reasoning>

<answer>
[Provide only the final answer here]
</answer>

Sampling Parameters: Temperature = 0.7, Samples = 10

```

Figure 7: Self-consistency prompt template (10 samples).

```

Self-Refine Prompt

Solve this problem step by step.
Question: {question}

Provide your response in the following format:

<reasoning>
[Show your step-by-step work and explain your thinking process here]
</reasoning>

<answer>
[Provide only the final answer here]
</answer>

Sampling Parameters: Temperature = 0, Samples = 1

```

Figure 8: Self-refine prompt template.

## F Computational Cost Analysis

We compare the computational costs of MARS against recursive self-improvement baselines: Gödel Agent (Yin et al., 2025) and Meta Agent Search (ADAS) (Hu et al., 2025a). Costs are estimated based on reported experimental configurations and current API pricing.

### F.1 Baseline Method Costs

**Meta Agent Search.** As reported in (Hu et al., 2025a), Meta Agent Search runs for **25 iterations**, where each iteration involves: (1) the meta agent (GPT-4) programming a new agent design, (2) self-reflection refinement (2 iterations per proposal, plus up to 3 error-correction refinements), and (3) evaluation on validation data using GPT-3.5. The authors report a total cost of approximately **\$300** across four benchmarks (DROP, MGSM, MMLU, GPQA). This high cost stems from the extensive GPT-4 usage for iterative agent design and the growing archive of discovered agents that must be included in each subsequent prompt.

**Gödel Agent.** As reported in (Yin et al., 2025), Gödel Agent performs **30 recursive self-improvements** across four benchmarks, with a total cost of approximately **\$15**. The framework uses GPT-4o for self-modification and GPT-3.5 for policy evaluation. The reduced cost compared to Meta Agent Search is attributed to continuous self-optimization that enables faster convergence. However, the authors note that the main cost driver is the continuously growing historical memory, suggesting that longer optimization runs would incur substantially higher costs.

### F.2 MARS Cost Estimation

MARS operates in a **single recurrence cycle** with three phases:

- Evaluation Phase:** The analyzer model processes each failed question to produce structured diagnoses. For a typical benchmark with  $\sim 200$  failed questions, this requires  $\sim 200$  API calls.
- Failure Allocation Phase:** Pure computa-

Enhancement Analyzer Prompt

Analyze this failed GPQA (Graduate-Level Google-Proof Q&A) question. The model used "{strategy}" prompting strategy.

Domain: {domain}  
 Subdomain: {subdomain}

Question: {question}

Answer Choices:  
 {choices\_text}

Correct Answer: {correct\_letter} - {correct\_answer}  
 Model's Answer: {predicted\_letter} - {model\_answer}

Model's Reasoning (excerpt): {reasoning}

Provide a comprehensive analysis in JSON format:

```

{
  "question_type": "<factual/conceptual/calculation/application/analysis/comparison>",
  "topics": ["<specific scientific topic 1>", "<topic 2>"],
  "error_type": "<conceptual_misunderstanding/calculation_error/misreading/...>",
  "root_cause": "<detailed explanation of why wrong answer was chosen>",
  "specific_mistake": "<exact reasoning step that failed>",
  "requires_knowledge": ["<scientific knowledge 1>", "<knowledge 2>"],
  "difficulty_factors": ["<what makes this challenging for AI>"]
}

```

Focus on why the model chose the wrong answer in this multiple-choice context.

Figure 9: Enhancement Analyzer prompt template for individual failure analysis.

Method	Baseline	Concise	Reasoning	R+C	Hybrid
Zero-shot	11.82	12.73	11.82	13.64	19.10
Zero-shot-CoT	16.40	18.18	15.45	13.64	22.70
Self-refine	36.36	37.27	39.09	37.27	48.20
Self-consistency	18.20	19.09	17.27	22.73	32.60

Table 11: GPQA performance (%) with Qwen-generated enhancements.

934 tional grouping with no API calls required.

935 3. **Enhancement Generation Phase:** The syn-  
 936 thesizer generates enhancements for each  
 937 type-topic group. With ~15–20 groups per  
 938 benchmark, this requires ~40–60 API calls  
 939 (including concise, reasoning, and specific  
 940 variants).

941 Using GPT-3.5-turbo (\$0.0005/1K input,  
 942 \$0.0015/1K output tokens) for both analysis and  
 943 synthesis, the estimated cost per benchmark is:

### 944 F.3 Cost Comparison Summary

945 Table 14 summarizes the computational require-  
 946 ments across methods.

### F.4 Analysis

The cost reduction achieved by MARS stems from three design decisions:

- 950 1. **Single-cycle learning:** While recursive meth-  
 951 ods require 25–30 iterations to converge,  
 952 MARS consolidates learning into one pass,  
 953 eliminating the multiplicative cost of iteration.
- 954 2. **No growing context:** Recursive methods ac-  
 955 cumulate historical memory (Gödel Agent) or  
 956 an archive of discovered agents (Meta Agent  
 957 Search), causing token usage to grow with  
 958 each iteration. MARS processes fixed-size  
 959 inputs throughout.
- 960 3. **Efficient model selection:** MARS uses GPT-  
 961 3.5 for both analysis and synthesis, while re-

**Enhancement Synthesizer Prompt**

Analyze these {num\_failures} failed GPQA questions that share:

Domain: {domain}  
 Question Type: {question\_type}  
 Topics: {topics}  
 Prompting Strategy Used: {strategy}

Sample Failures:  
 {failures\_summary}

Common Error Patterns: {error\_patterns}  
 Required Knowledge: {required\_knowledge}  
 Difficulty Factors: {difficulty\_factors}

Create a targeted enhancement strategy in JSON format:

```
{
  "common_mistakes": ["<specific mistake pattern 1>", "<pattern 2>"],
  "key_warnings": ["<critical warning for this topic 1>", "<warning 2>"],
  "verification_steps": ["<verification step 1>", "<step 2>"],
  "topic_specific_guidance": "<detailed guidance for these topics>",
  "type_specific_approach": "<approach for this question type>",
  "enhanced_prompt_addition": "<concise 2-3 sentence prompt addition>"
}
```

Make it specific to {question\_type} questions about {topics} in {domain}.

Figure 10: Enhancement Synthesizer prompt template for pattern-based enhancement generation.

Method	Baseline	Concise	Reasoning	R+C	Hybrid
Zero-shot	11.82	12.73	11.82	13.64	19.10
Zero-shot-CoT	16.40	18.18	15.45	13.64	22.70
Self-refine	36.36	37.27	39.09	37.27	48.20
Self-consistency	18.20	19.09	17.27	22.73	32.60

Table 12: OMNI-math performance (%) with Qwen-generated enhancements.

Phase	API Calls	Tokens (K)	Est. Cost
Evaluation	~200	~400	~\$0.40
Enhancement Gen.	~50	~150	~\$0.15
<b>Total (per benchmark)</b>	~250	~550	~\$0.55
<b>Total (4 benchmarks)</b>	~1,000	~2,200	~\$2.20

Table 13: Estimated MARS computational cost using GPT-3.5-turbo.

cursive methods require GPT-4/GPT-4o for meta-level reasoning. As shown in Appendix D, enhancement quality is robust to generator model choice.

The cost-performance trade-off favors MARS: at 136× lower cost than Meta Agent Search and 6.8× lower than Gödel Agent, MARS achieves comparable or superior performance (Table 4). For applica-

tions requiring maximum performance regardless of cost, MARS can be applied iteratively—using enhanced prompts to generate new failure sets for further refinement—while still maintaining substantial cost advantages over recursive baselines.

## G MARS Implementation Details

This appendix provides key code snippets from the MARS implementation.

### G.1 Data Structures

Listing 1 defines the core data structures corresponding to Equations 1–3.

```
1 @dataclass
2 class IndividualFailureAnalysis:
3     """Structured analysis A_i = (tau_i, T_i, epsilon_i, rho_i,
4     mu_i)"""
5     question_id: str
6     question_text: str
```

## Concise Enhancement Sample

```

Answer the following question.
Question: {question}

## SPECIALIZED GUIDANCE FOR ALGEBRA EQUATIONS (ZERO_SHOT)
### Critical Warnings by Question Type:

**Calculation questions about algebra/number theory** (7 failures):
[!] Ensure that all algebraic relationships and constraints are clearly
defined before attempting calculations.
[!] Be cautious of jumping to conclusions without verifying all potential
solutions and constraints.
-> When solving calculation problems in algebra and number theory,
carefully define all variables and relationships before proceeding.
Consider all possible integer solutions and verify constraints.

**Calculation questions about algebra/equations** (5 failures):
[!] Be cautious of prematurely concluding an answer without completing
all necessary steps.
[!] Ensure all potential simplifications and substitutions are considered.
-> Approach this problem by identifying any potential simplifications
or substitutions. Carefully verify each algebraic manipulation step.

**Analysis questions about functional equations** (3 failures):
[!] Ensure all solutions satisfy the functional equation for all x and y.
-> When solving functional equations, verify each proposed solution
thoroughly. Explore transformations like symmetry or substitution.
...

```

Enhancement Type: Warning Indicators [!] + Action Sequences (->)

Figure 11: Concise enhancement template structure. Warnings identify failure patterns; arrows provide actionable guidance.

Method	Iterations	Meta Model	Eval Model	Cost (4 benchmarks)	Relative
Meta Agent Search	25	GPT-4	GPT-3.5	~\$300	136×
Gödel Agent	30	GPT-4o	GPT-3.5	~\$15	6.8×
<b>MARS (Ours)</b>	<b>1</b>	GPT-3.5	GPT-3.5	~\$2.20	<b>1×</b>

Table 14: Computational cost comparison across self-improvement methods. Costs based on reported values (Yin et al., 2025; Hu et al., 2025a) and our estimates.

```

988 6 question_type: str # tau_i in Y (question type)
989 7 topics: List[str] # T_i subset of D (topic set)
990 8 error_type: str # epsilon_i in E (error taxonomy)
991 9 root_cause: str # rho_i (reasoning deficit)
992 10 specific_mistake: str # mu_i (divergence point)
993 11 requires_knowledge: List[str]
994 12 difficulty_factors: List[str]
995 13
996 14 @dataclass
997 15 class QuestionTypeTopicGroup:
998 16     """Group G_j with error profile Psi_j = (E_j, R_j, F_j)"""
999 17     question_type: str
1000 18     topics: List[str]
1001 19     failures: List[IndividualFailureAnalysis]
1002 20     common_error_patterns: List[str] # E_j
1003 21     shared_root_causes: List[str] # R_j
1004 22     required_knowledge: Set[str]
1005 23     key_difficulty_factors: List[str] # F_j
1006 24
1007 25 @dataclass
1008 26 class TypeTopicEnhancement:
1009 27     """Enhancement variants for type-topic group:
1010 28     E^c(c): concise, E^r(r): reasoning, E^c(c+r): specific"""
1011 29     question_type: str
1012 30     topics: List[str]
1013 31     num_questions: int # |G_j| for weight w_j
1014 32     # For concise enhancement E^c(c):
1015 33     key_warnings: List[str] # Warning indicators [!]
1016 34     # For reasoning enhancement E^r(r):
1017 35     enhanced_prompt_addition: str # Process-oriented

```

```

1018 guidance
1019 # For specific enhancement E^c(c+r) = concise + reasoning:
1020 common_mistakes: List[str] # Explicit error patterns
1021 (x)
1022 verification_steps: List[str] # Validation actions (+)
1023 type_specific_approach: str # Methodological guidance

```

Listing 1: Core data structures for MARS pipeline.

## G.2 Phase 1: Individual Failure Analysis 1025

Listing 2 shows the evaluation phase that produces structured analyses  $\mathcal{A}_i$  for each failed question. 1026

```

1027 1028
1029 1 def analyze_individual_failure(self, failure: Dict,
1030 2 strategy: str) ->
1031 3     IndividualFailureAnalysis:
1032 4     question = failure.get('question', '')
1033 5     correct_answer = failure.get('correct_answer', '')
1034 6     model_answer = failure.get('predicted_answer', '')
1035 7
1036 8     prompt = f"""Analyze this failed question using "{strategy
1037 9 }" strategy.
1038 10 Question: {question[:2000]}
1039 11 Correct: {correct_answer[:500]}
1040 12 Model Answer: {model_answer[:500]}
1041

```

## Reasoning Enhancement Sample

Answer the following question.  
Question: {question}

## SPECIALIZED GUIDANCE FOR ALGEBRA EQUATIONS (ZERO\_SHOT)  
### Key Considerations by Problem Type:

- \* Calculation (algebra/number theory): When solving calculation problems in algebra and number theory, carefully define all variables and relationships before proceeding. Consider all possible integer solutions and verify that they satisfy the given constraints.
- \* Calculation (algebra/equations): Approach this problem by identifying any potential simplifications or substitutions that could ease calculations. Carefully verify each algebraic manipulation step.
- \* Analysis (functional equations/real analysis): When solving functional equations, verify each proposed solution thoroughly. Consider specific values of  $x$  and  $y$ , and explore transformations like symmetry.
- \* Reasoning (functional equations/real analysis): When tackling functional equations, consider starting with simple functions like linear or constant. Use substitutions to simplify the equation.
- \* Calculation (algebra/systems of equations): When solving algebraic systems of equations, first seek patterns or symmetries that might simplify the problem. Verify solutions by substituting them back.
- \* Calculation (algebra/optimization): When tackling algebraic optimization problems, carefully analyze the constraints and consider algebraic manipulations that simplify the expression.

*Enhancement Type: Process-Oriented Guidance (\*)*

Figure 12: Reasoning enhancement template structure. Bullet points provide process-oriented guidance for each problem type.

## Reasoning + Concise Enhancement Sample

Answer the following question.

Question: {question}

## SPECIALIZED GUIDANCE FOR ALGEBRA EQUATIONS (ZERO\_SHOT)

### Detailed Guidance by Question Type and Topic:

**\*\*Calculation - algebra & number theory\*\*** (7 failures):

Common Mistakes:

- x Failing to correctly establish algebraic relationships from word problems, leading to incorrect equation setups.
- x Misapplying integer constraints and failing to consider all possible integer solutions that satisfy the given conditions.
- x Jumping to conclusions without systematically verifying each candidate solution against all problem constraints.

Verification Steps:

- + Break down the problem into smaller, manageable parts and clearly define all variables and their relationships before solving.
- + Verify each step of the calculation, especially when dealing with integer constraints, by checking edge cases.
- + Ensure all solutions satisfy the original equations by substituting back and confirming consistency with problem conditions.

Approach: In calculation questions involving algebra and number theory, start by translating the problem into mathematical expressions. Use algebraic manipulation to simplify and solve these equations, and apply number theory principles to handle integer constraints. Consider multiple approaches such as substitution, factorization, or case analysis to ensure comprehensive exploration of solutions.

**\*\*Calculation - algebra & equations\*\*** (5 failures):

Common Mistakes:

- x Failure to complete calculations, leading to incomplete answers.
- x Incorrect application of algebraic manipulations or simplifications.

Verification Steps:

- + Review the problem for potential simplifications or substitutions.
- + Verify each step of the algebraic manipulation to ensure no errors.
- + Cross-check final answer by substituting back into original equations.

Approach: For calculation questions in algebra, emphasize a step-by-step approach that includes checking for simplifications and verifying each manipulation. Leverage known algebraic identities or properties.

...

*Enhancement Type: Common Mistakes (x) + Verification Steps (+) + Approach*

Figure 13: Specific enhancement template structure. Combines explicit mistake warnings, verification steps, and methodological guidance.

```

1042 Provide JSON analysis:
1043 {{
1044     "question_type": "<factual/conceptual/calculation/
1045     application>",
1046     "topics": ["<topic_1>", "<topic_2>"],
1047     "error_type": "<conceptual_misunderstanding/
1048     calculation_error/...>",
1049     "root_cause": "<fundamental reasoning deficit>",
1050     "specific_mistake": "<exact step where logic diverged>",
1051     "requires_knowledge": ["<knowledge_1>"],
1052     "difficulty_factors": ["<factor_1>"]
1053 }}"""
1054
1055 result = call_llm(self.client, self.model,
1056                 [{"role": "user", "content": prompt}],
1057                 temperature=0.3, max_tokens=800)
1058 data = json.loads(extract_json(result))
1059
1060 return IndividualFailureAnalysis(
1061     question_type=data.get('question_type'),
1062     topics=data.get('topics', []),
1063     error_type=data.get('error_type'),
1064     root_cause=data.get('root_cause', ''),
1065     specific_mistake=data.get('specific_mistake', ''),
1066     requires_knowledge=data.get('requires_knowledge', []),
1067     difficulty_factors=data.get('difficulty_factors', []))

```

Listing 2: Code for individual failure analysis (Evaluation phase).

### G.3 Phase 2: Type-Topic Grouping

Listing 3 implements the grouping function  $\kappa$  (Equation 2) that partitions analyses into groups  $\mathcal{G} = \{G_j\}$ .

```

1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
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

```

```

1 def group_by_type_topic(self, analyses: List[
2     IndividualFailureAnalysis]
3     ) -> List[QuestionTypeTopicGroup]:
4     """Apply grouping function kappa: A -> Y x 2^D"""
5
6     # Partition by composite key (tau_i, T_i)
7     groups = defaultdict(list)
8     for analysis in analyses:
9         key = (analysis.question_type,
10              frozenset(analysis.topics[:2]))
11         groups[key].append(analysis)
12
13     type_topic_groups = []
14     for (q_type, topics), group_analyses in groups.items():
15         # Aggregate into error profile Psi_j
16         error_patterns = [a.error_type for a in group_analyses]
17
18         root_causes = [a.root_cause for a in group_analyses]
19         required_knowledge = set()
20         difficulty_factors = []
21
22         for a in group_analyses:
23             required_knowledge.update(a.requires_knowledge)
24             difficulty_factors.extend(a.difficulty_factors)
25
26         group = QuestionTypeTopicGroup(
27             question_type=q_type,
28             topics=list(topics),
29             failures=group_analyses,
30             common_error_patterns=list(set(error_patterns)),
31             shared_root_causes=list(set(root_causes)),
32             required_knowledge=required_knowledge,
33             key_difficulty_factors=list(set(difficulty_factors))
34         ))
35         type_topic_groups.append(group)
36
37     return type_topic_groups

```

Listing 3: Code for failure allocation via type-topic grouping.

### G.4 Phase 3: Enhancement Generation

Listing 4 implements the enhancement function  $\xi$  (Equation 4) and prompt aggregation (Equation 5).

The three enhancement variants are: (1) **concise**  $E^{(c)}$ : warning indicators + action sequences; (2) **reasoning**  $E^{(r)}$ : process-oriented guidance; and (3) **specific**  $E^{(c+r)}$ : the combination of concise and reasoning, providing common mistakes, verification steps, and methodological approach.

```

1 def create_enhanced_prompts(self, base_prompt: str,
2                             enhancements: List[
3                                 TypeTopicEnhancement],
4                             strategy: str, category: str) ->
5     Dict[str, str]:
6     """Generate P^(c) and P^(r) via weighted aggregation"""
7
8     # Sort by |G_j| for weight w_j (Eq. 5)
9     sorted_enh = sorted(enhancements,
10                        key=lambda e: e.num_questions, reverse
11                        =True)
12     all_prompts = {}
13
14     # Concise enhancement E^(c): warnings + action sequences
15     if 'concise' in self.enhancement_types:
16         text = f"\n## GUIDANCE FOR {category.upper()}\n"
17         text += "### Critical Warnings by Question Type:\n\n"
18         for enh in sorted_enh[:8]: # Top-weighted groups
19             text += f"*{enh.question_type} ({'/'}.join(enh.
20                 topics))* "
21             text += f"({enh.num_questions} failures):\n"
22             text += "(!" + " | ".join(enh.key_warnings[:3])
23             + "\n"
24             text += f" -> {enh.enhanced_prompt_addition}\n\n"
25         all_prompts['concise'] = base_prompt + text # P' = P
26         + E
27
28     # Reasoning enhancement E^(r): process-oriented hints
29     if 'reasoning' in self.enhancement_types:
30         text = f"\n## GUIDANCE FOR {category.upper()}\n"
31         text += "### Key Considerations by Problem Type:\n\n"
32         for enh in sorted_enh[:6]:
33             text += f"*{enh.question_type} ({'/'}.join(enh.
34                 topics))* "
35             text += f"({enh.enhanced_prompt_addition}\n"
36             all_prompts['reasoning'] = base_prompt + text
37
38     # Specific enhancement E^(c+r): combines concise +
39     # reasoning
40     # Includes: mistakes (concise) + verification (concise) +
41     # approach (reasoning)
42     if 'specific' in self.enhancement_types:
43         text = f"\n## GUIDANCE FOR {category.upper()}\n"
44         for enh in sorted_enh[:10]:
45             text += f"*{enh.question_type} - {'& '}.join(enh.
46                 topics))*\n"
47             # From concise: explicit warnings as mistake
48             # patterns
49             text += "Common Mistakes:\n"
50             for m in enh.common_mistakes[:3]: text += f" x {m}
51             }\n"
52             # From concise: action sequences as verification
53             text += "Verification Steps:\n"
54             for s in enh.verification_steps[:4]: text += f" +
55             {s}\n"
56             # From reasoning: process-oriented approach
57             text += f"Approach: {enh.type_specific_approach}\n"
58             all_prompts['specific'] = base_prompt + text
59
60     return all_prompts

```

Listing 4: Code for enhancement generation and prompt aggregation.

### G.5 Hybrid Selection

Listing 5 implements the hybrid strategy (Equation 6) that selects optimal enhancement per category.

```

1 def select_hybrid_enhancement(self, val_data: Dict[str, List],

```

```

1190 2          enhancements: Dict) -> Dict[str,
1191 3      str]:
1192 4      """Select E*_c = argmax Acc(E, V_c) for each category c
1193 5      where E in {E*(c), E*(r), E*(c+r)} (concise, reasoning,
1194 6      specific)"""
1195 7
1196 8      optimal = {}
1197 9      # specific = concise + reasoning (c+r)
1198 10     etypes = ['concise', 'reasoning', 'specific']
1199 11
1200 12     for category, val_questions in val_data.items():
1201 13         best_acc, best_type = 0, 'concise'
1202 14
1203 15         for etype in etypes:
1204 16             enhanced_prompt = enhancements.get(f"{category}_{
1205 17             etype}")
1206 18             if not enhanced_prompt: continue
1207 19
1208 20             correct = sum(1 for q in val_questions
1209 21                 if self.evaluate(enhanced_prompt, q)
1210 22                 == q['correct_answer'])
1211 23             accuracy = correct / len(val_questions)
1212 24
1213 25             if accuracy > best_acc:
1214 26                 best_acc, best_type = accuracy, etype
1215 27
1216 28             optimal[category] = best_type
1217 29             print(f" {category}: '{best_type}' (acc: {best_acc
1218 30             :.1%})")
1219
1220 31     return optimal

```

Listing 5: Code for hybrid enhancement selection.