

# Atlas: Orchestrating Heterogeneous Models and Tools for Multi-Domain Complex Reasoning

Jinyang Wu<sup>1\*†</sup>, Guocheng Zhai<sup>1\*</sup>, Ruihan Jin<sup>1\*</sup>, Jiahao Yuan<sup>3</sup>, Yuhao Shen<sup>2</sup>,  
Shuai Zhang<sup>1</sup>, Zhengqi Wen<sup>1</sup>, Jianhua Tao<sup>1‡</sup>

<sup>1</sup>Tsinghua University, <sup>2</sup>Zhejiang University,

<sup>3</sup>East China Normal University

wu-jy23@mails.tsinghua.edu.cn

## Abstract

The integration of large language models (LLMs) with external tools has significantly expanded the capabilities of AI agents. However, as the diversity of both LLMs and tools increases, selecting the optimal model-tool combination becomes a high-dimensional optimization challenge. Existing approaches often rely on a single model or fixed tool-calling logic, failing to exploit the performance variations across heterogeneous model-tool pairs. In this paper, we present ATLAS (Adaptive Tool-LLM Alignment and Synergistic Invocation), a dual-path framework for dynamic tool usage in cross-domain complex reasoning. ATLAS operates via a dual-path approach: (1) **training-free cluster-based routing** that exploits empirical priors for domain-specific alignment, and (2) **RL-based multi-step routing** that explores autonomous trajectories for out-of-distribution generalization. Extensive experiments across 15 benchmarks demonstrate that our method outperforms closed-source models like GPT-4o, surpassing existing routing methods on both in-distribution (+10.1%) and out-of-distribution (+13.1%) tasks. Furthermore, our framework shows significant gains in visual reasoning by orchestrating specialized multi-modal tools.

## 1 Introduction

Large language models (LLMs) have evolved from static problem solvers into collaborative reasoning engines through adaptive integration with external tools. These tools range from symbolic reasoning modules (Feng et al., 2025a) to real-time information retrieval APIs (Ma et al., 2025), significantly extending LLMs’ operational capabilities. As this LLM-tool ecosystem evolves, the synergy from multiple candidates increasingly surpasses the potential of either routing in model swarms (Yue et al.,

\* Equal Contribution.

† Project Lead.

‡ Corresponding Authors.

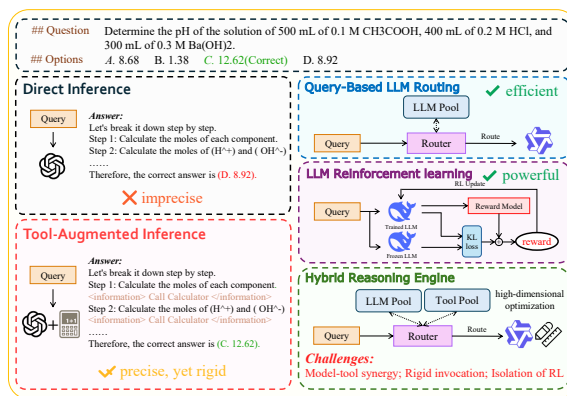


Figure 1: **Comparison of different LLM inference paradigms.** While routing (efficiency) and RL (performance optimization) present a promising approach, dynamic tool usage still faces significant challenges.

2025b) or tool augmentation (Dong et al., 2025) alone, highlighting the critical need for identifying the optimal model-tool combination.

Recent advances have focused on different aspects of these reasoning engines separately. **For tool usage**, existing frameworks (Kong et al., 2024; Wu et al., 2024) improve performance through task planning, yet relying on fixed logic that cannot dynamically adapt to different model capabilities or task requirements. **For LLM routing**, methods like ZOOTER (Lu et al., 2024) and RouterDC (Chen et al., 2024) optimize model selection through reward-guided learning and dual contrastive learning. Likewise, frameworks such as HybridLLM (Ding et al., 2024) and RouteLLM (Ong et al., 2024) combine strong and weak models for cost efficiency. However, these routing methods treat models as isolated execution units and fail to incorporate external tools, which could significantly enhance task performance. **For reinforcement learning (RL)**, methods such as RLHF (Ouyang et al., 2022) and PPO (Schulman et al., 2017) are explored to optimize reasoning capabilities in LLMs. RLAIIF (Lee et al., 2023) and DPO (Rafailov et al., 2023) bypass explicit reward modeling, stream-

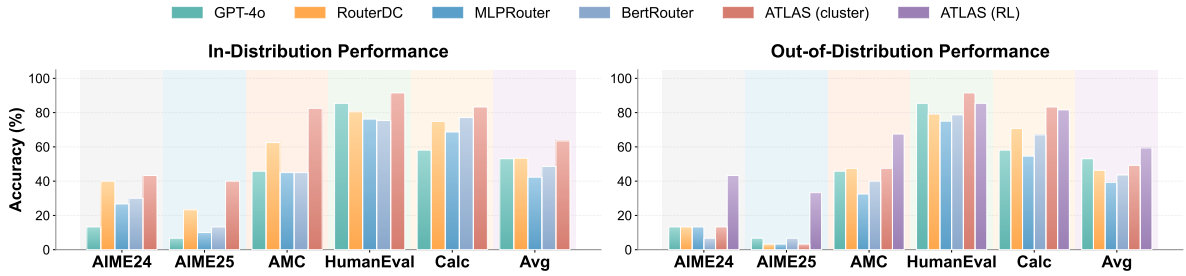


Figure 2: **Performance comparison on in-distribution and out-of-distribution settings.** Our ATLAS method consistently outperforms all baselines across diverse datasets, demonstrating superior generalization capability.

lining preference learning. Additionally, RouterR1 (Zhang et al., 2025a) allows models to deliberate internally before invoking auxiliary models. Recent works (Chen et al., 2025b; Jin et al., 2025a; Feng et al., 2025a) apply RL to tool usage, but miss the opportunity to integrate both models and tools to fully harness their combined strengths.

As shown in Figure 1, existing methods neglect the dynamic interplay of tool usage, LLM routing and RL, thus falling short especially when faced with the emerging diversity of LLMs and tools. This fundamental limitation manifests in three key challenges: (1) **Failure to leverage model-tool synergies:** LLM routing methods focus solely on model selection without integrating external tools, limiting their potential to enhance task performance; (2) **Rigid invocation and limited flexibility:** Existing tool usage methods rely on fixed, pre-configured invocation logic that hinders adaptability and scalability, preventing reasoning engines from dynamically optimizing model-tool combinations in open-domain tasks; (3) **Isolated optimization of RL:** Even advanced RL approaches focus on optimizing individual components in isolation, missing opportunities to jointly leverage model-tool synergies for complex reasoning.

To address these challenges, we propose ATLAS (Adaptive Tool-LLM Alignment and Synergistic Invocation), a generalizable framework that dynamically orchestrates optimal model-tool combinations. Our approach employs a dual-path approach to bridge the gap between empirical knowledge and open-domain reasoning. We firstly introduce **training-free cluster-based routing** that efficiently selects model-tool pairs by leveraging domain-specific expertise within a semantic embedding space. This approach exploits historical performance patterns for rapid, accurate routing in familiar domains. For generalized scenarios where explicit priors are absent, we utilize **RL-based multi-step routing** that iteratively

explores the model-tool combinations for superior execution paths. This bifurcated design effectively resolves the scalability challenges in high-dimensional search spaces while ensuring robustness. We conduct experiments on 15 benchmarks to evaluate the proposed ATLAS in both in-distribution and out-of-distribution settings. Empirical results shown in Figure 2 reveal that ATLAS achieves a superior performance across diverse tasks, which demonstrates its effectiveness as a new paradigm for tool-augmented reasoning agents. Our primary contributions are as follows:

- We introduce **ATLAS**, a generalizable agentic framework that explicitly optimizes heterogeneous synergies between diverse LLMs and tools, enabling dynamic and adaptive tool invocation for complex reasoning tasks.
- We propose a dual-path design that handles both domain-specific and open-domain tasks: (1) **training-free cluster-based routing** for efficient selection using domain expertise, and (2) **RL-driven multi-step routing** for generalizing across unfamiliar tasks via iterative exploration.
- Experiments across 9 tasks and 15 benchmarks show that ATLAS outperforms top-performing closed-source LLMs and powerful routing methods on multi-domain tasks and exhibits robust adaptability in multi-modal scenarios.

## 2 Related Work

**Query-based LLM Routing.** As the landscape of LLMs continues to evolve, query-based routing has become crucial in reasoning engines for balancing performance and computational efficiency by dynamically selecting the most appropriate model for each query. Early approaches rely on reward-guided (Lu et al., 2024) and contrastive learning strategies (Chen et al., 2024) to improve routing

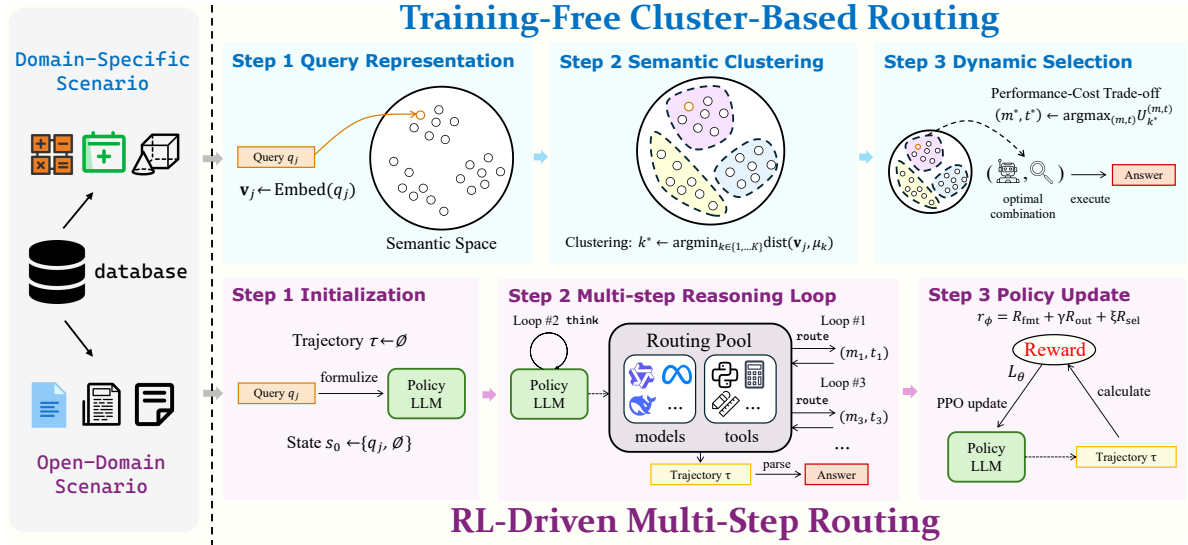


Figure 3: Overview of Adaptive Tool-LLM Alignment and Synergistic Invocation (ATLAS). The framework operates via a dual-path approach: (1) Training-free Cluster-based Routing; and (2) RL-driven Multi-step Routing.

accuracy. Existing methods balance computational cost with performance through query-level orchestration (Ding et al., 2024; Ong et al., 2024; Zhang et al., 2025b), model cascading (Chen et al., 2023), adaptive selection (Feng et al., 2025b; Wang et al., 2025; Jin et al., 2025b), and budget allocation (Mei et al., 2025). Further, the integration of routing within reasoning frameworks (Yue et al., 2025b; Pan et al., 2025) enhances the performance boundaries. However, these existing methods often treat LLMs as isolated execution units, neglecting the synergies between specific model capabilities and external tool interfaces. Our approach addresses this gap by jointly optimizing model-tool combinations, enabling a more adaptive, scalable, and effective reasoning engine capable of dynamically integrating the strengths of both models and tools.

**Reinforcement Learning for LLM.** Reinforcement learning (RL) has been widely applied to optimize LLMs for aligning with complex human preferences and improving reasoning tasks. The paradigm has evolved from reward-model-based approaches like RLHF (Ouyang et al., 2022) and PPO (Schulman et al., 2017) to more efficient frameworks such as DPO (Rafailov et al., 2023), which bypass explicit reward modeling to streamline preference learning. RL has also been applied to optimize routing decisions, with approaches like Router-R1 (Zhang et al., 2025a) allowing models to deliberate internally before invoking auxiliary models. Recent works (Chen et al., 2025b; Jin et al., 2025a; Feng et al., 2025a) investigate the application of RL for tool usage. While these methods

reveal the potential of RL in optimizing reasoning trajectories, they primarily focus on single-model or single-tool optimization, overlooking the large potential for combined synergies. ATLAS extends this by employing RL to jointly optimize model-tool combinations, enabling more adaptive and efficient reasoning.

### 3 Methodology

We present the ATLAS framework (Figure 3), which combines a two-tier strategy: *Training-free Cluster-based Routing* (§ 3.1) to enable quick decision-making; and *RL-driven Multi-step Routing* (§ 3.2) to handle more complex open-domain tasks that require iterative model-tool interactions.

#### 3.1 Training-Free Cluster-Based Routing

We hypothesize that the optimal model-tool combination is query-dependent and exhibits semantic locality. Consequently, the empirical strategy approximates the optimal routing function  $f : \mathcal{Q} \rightarrow \mathcal{S}$  by leveraging historical metadata. We define the search space as the Cartesian product  $\mathcal{S} = \mathcal{M} \times \mathcal{T}$ , where  $\mathcal{M} = \{m_1, \dots, m_M\}$  denotes the set of candidate LLMs and  $\mathcal{T} = \{t_1, \dots, t_T\}$  represents the available tools.

Given  $\mathcal{Q}_{\text{train}} = \{q_i\}_{i=1}^N$  denote the training queries set, we map each query  $q_i$  into a  $D$ -dimensional latent manifold using a pre-trained encoder:  $\mathbf{v}_i = \mathcal{E}(q_i) \in \mathbb{R}^D$ . To capture the semantic task distribution, we partition the embedding space into  $K$  disjoint clusters  $\{\mathcal{C}_k\}_{k=1}^K$  by minimiz-

ing the inertia:

$$\min_{\{\mu_k\}_{k=1}^K} \sum_{k=1}^K \sum_{\mathbf{v}_i \in \mathcal{C}_k} \|\mathbf{v}_i - \mu_k\|^2, \quad (1)$$

where  $\mu_k$  represents the semantic centroid of cluster  $\mathcal{C}_k$ . The clustering process effectively groups queries with similar reasoning requirements and tool affinities.

We derive empirical statistics from the training observations for each modal-tool pair  $(m, t) \in \mathcal{S}$  within the cluster  $\mathcal{C}_k$ . The empirical accuracy is defined as the success rate of the pair  $(m, t)$  on the cluster  $\mathcal{C}_k$ :

$$\widehat{\text{Acc}}_k^{(m,t)} = \frac{1}{|\mathcal{C}_k|} \sum_{q_i \in \mathcal{C}_k} \mathbb{1}[(m, t) \text{ solves } q_i], \quad (2)$$

where  $\mathbb{1}[\cdot]$  denotes the indicator function.

Simultaneously, we model the operational cost to account for resource consumption, which is computed based on the average token throughput observed during the profiling phase:

$$\widehat{\text{Cost}}_k^{(m,t)} = \bar{N}_{\text{in}}^{(m,t)} \cdot P_{\text{in}}^{(m,t)} + \bar{N}_{\text{out}}^{(m,t)} \cdot P_{\text{out}}^{(m,t)}, \quad (3)$$

where  $\bar{N}_{\text{in}}$  and  $\bar{N}_{\text{out}}$  represent the mean input and output token counts for the cluster, while  $P_{\text{in}}$  and  $P_{\text{out}}$  denote their respective unit prices.

To facilitate a flexible trade-off between reasoning performance and inference cost, we define a cluster-level utility score  $\mathcal{U}_k(m, t)$  as:

$$\mathcal{U}_k(m, t) = (1 - \alpha) \cdot \widehat{\text{Acc}}_k^{(m,t)} - \alpha \cdot \widehat{\text{Cost}}_k^{(m,t)}, \quad (4)$$

where  $\alpha \in [0, 1]$  is a hyperparameter that balances the performance-cost trade-off.

At inference time, the framework performs low-latency orchestration by projecting a novel query  $q_j$  into the latent manifold  $\mathbf{v}_j = \mathcal{E}(q_j)$ . The routing is executed via a proximal cluster lookup, where the query is assigned to  $k^* = \arg \min_k \|\mathbf{v}_j - \mu_k\|$ . Subsequently, the system retrieves the optimal model-tool pair for execution:

$$(m^*, t^*) = \arg \max_{(m,t) \in \mathcal{S}} \mathcal{U}_{k^*}(m, t). \quad (5)$$

By caching heterogeneous synergies within the embedding space, this empirical strategy enables real-time, cost-aware tool invocation with constant-time complexity relative to the number of clusters.

### 3.2 RL-Driven Multi-Step Routing

While the empirical strategy excels in low-latency routing, it is inherently limited by its reliance on a single-shot decision. To address complex tasks that demand multi-round reasoning and iterative model-tool interactions, we introduce an RL-driven strategy that instantiates the router as an autonomous agent capable of interleaving internal reasoning with external invocation.

We model this process as a sequential decision task over a maximum horizon  $T_{\text{max}}$ . For a given query  $q_j$ , the agent maintains an evolving state  $s_t = \{q_j, C_t\}$ , where  $C_t$  represents the accumulated context of previous reasoning trajectories and tool outputs. At each step  $t$ , the policy  $\pi_\theta$  samples an action  $a_t$  from the augmented action space  $\mathcal{A}$ , comprising two types of operations: (1) *Internal Reasoning* (think), where the agent performs local chain-of-thought processing to decompose complex queries or synthesize intermediate results; and (2) *Dynamic Routing* (route), where the agent selects a specific model-tool pair  $(m, t) \in \mathcal{S}$  from the routing pool to gather external observations  $o_t$ . This iterative loop ensures that the agent can adaptively refine its search space based on real-time feedback from the environment until an answer is extracted or the maximum step limit is reached.

To optimize this decision-making process, we train the policy  $\pi$  using Proximal Policy Optimization (PPO) (Schulman et al., 2017), which maximizes the following regularized objective:

$$\max_{\pi} \mathbb{E}_{q \sim \mathcal{D}, \tau \sim \pi} \left[ r_\phi(q, \tau) - \beta \log \frac{\pi(\tau|q; \mathcal{P})}{\pi_{\text{ref}}(\tau|q; \mathcal{P})} \right], \quad (6)$$

where  $\tau$  is the interaction trajectory,  $\pi_{\text{ref}}$  is a reference policy to ensure training stability, and  $\beta$  is the KL-regularization coefficient.

We design the reward function  $r_\phi$  as a composite of three finely-tuned rule-based signals (detailed in Appendix A.2) including format reward ( $\mathcal{R}_{\text{fmt}}$ ), outcome reward ( $\mathcal{R}_{\text{out}}$ ) and model selection reward ( $\mathcal{R}_{\text{sel}}$ ), bridging the gap between structured execution and task correctness, formally:

- **Format Reward ( $\mathcal{R}_{\text{fmt}}$ ):** A signal enforces structural integrity by penalizing trajectories that deviate from the predefined format and tool-invocation syntax.
- **Outcome Reward ( $\mathcal{R}_{\text{out}}$ ):** A binary signal that directly aligns the policy with task correctness.

- **Model Selection Reward ( $\mathcal{R}_{\text{sel}}$ ):** A penalty-based signal guides the agent toward optimal efficiency by penalizing the selection of sub-optimal models.

The final reward is computed as:

$$r_\phi = \mathcal{R}_{\text{fnt}} + \gamma\mathcal{R}_{\text{out}} + \xi\mathcal{R}_{\text{sel}}, \quad (7)$$

where  $\gamma$  and  $\xi$  are hyperparameters. This framework facilitates autonomous orchestration, as the model learns to assess the sufficiency of its internal state before invoking external resources. By decoupling routing logic via the  $\mathcal{R}_{\text{sel}}$  signal, ATLAS internalizes the fundamental alignment between domains and tool utilization rather than memorizing rigid model-tool mappings. This design ensures that the routing policy captures the essential characteristics of expertise distribution, remaining robust and generalizable even as the available tools and models evolve in dynamic environments.

## 4 Experiments

This section presents a comprehensive evaluation of ATLAS, covering main results across multi-domain benchmarks (§ 4.2), multi-modal visual reasoning (§ 4.3), model-tool pool extensions (§ 4.4), and further analysis on reasoning boundaries, model-tool alignment preferences, and RL convergence dynamics (§ 4.5).

### 4.1 Experimental Settings

**Models Selection.** To evaluate ATLAS’s generalization across model architectures and scales, we select six heterogeneous open-source LLMs: Qwen2.5-7B-Instruct (Yang et al., 2024a), Llama-3.1-8B-Instruct (Dubey et al., 2024), InternLM3-8B-Instruct (Cai et al., 2024), DeepSeek-R1-Distill-Qwen-7B (Guo et al., 2025), Qwen2.5-Coder-7B-Instruct (Hui et al., 2024), and a multi-modal LLM Qwen3-8B-VL-Instruct (Yang et al., 2025). This diverse selection allows us to observe how different models synergize with specific external tools.

**Tool Definition.** We introduce two tool sets for textual and visual reasoning:

- **Foundation Tools:** This set includes four essential tools: (1) *Code Interpreter*, a Python execution environment for algorithmic and logical verification; (2) *Web Search* for retrieving real-time open-domain information; (3) *Calculator* for high-precision numerical

computation; and (4) *Process Reward Model (PRM)* for scoring and ranking model outputs.

- **Multi-modal Tools:** (1) *Qwen3-Chart* for chart data extraction; (2) *Qwen3-Counting* for enumerating objects in images; (3) *Qwen3-Geo* for parsing geometric properties and performing post-hoc self-verification of geometric proofs; and (4) *Hunyuan-OCR* (Team et al., 2025) for text extraction from images. The first three tools use Qwen3-8B-VL with task-specific prompts, due to the underperformance of most existing specialized tools.

**Benchmarks and Baselines.** We evaluate on multi-domain tasks: (1) mathematical reasoning: AIME2024 (MAA, 2024), AIME2025 (MAA, 2025), AMC (Lightman et al., 2023); (2) code generation: HumanEval (Chen, 2021), MBPP (Austin et al., 2021); (3) arithmetic reasoning: Calculator (Wu et al., 2025d); (4) commonsense reasoning: NQ (Kwiatkowski et al., 2019), WebQ (Berant et al., 2013); (5) logical reasoning: LogiQA2 (Liu et al., 2023); (6) scientific reasoning: GPQA (Rein et al., 2024). Furthermore, we extend our evaluations to multi-modal benchmarks, including ChartQA (Masry et al., 2022), Geometry3K (Lu et al., 2021), TallyQA (Acharya et al., 2019), CountBench (Paiss et al., 2023), and TableVQA (Kim et al., 2024). We use accuracy as the primary metric. Baselines include **Zero-shot/Few-shot Router**, **Random Router**, **RouterDC** (Chen et al., 2024), **MLPRouter** (Hu et al., 2024), and **BertRouter** (Ong et al., 2024). Details are provided in Appendix B.

**Implementation Details.** For RL experiments, we adopt Qwen2.5-3B-Instruct as the policy model for model-tool selection. The policy is optimized with a batch size of 32 for 250 training steps, and the learning rate is set to  $1 \times 10^{-6}$ . More details are provided in Appendix B.4.

### 4.2 Main Results

Table 1 presents a comprehensive evaluation of our framework against various routing baselines across in-distribution and out-of-distribution tasks.

#### 4.2.1 In-Distribution Performance

Under the in-distribution setting, where training data for all tasks is accessible, ATLAS(cluster) achieves 63.5% average accuracy, surpassing the strongest baseline RouterDC by 10.1%. This advantage is pronounced on rigorous mathematical

Table 1: **Performance comparison across diverse tasks and domains.** *In-Distribution:* All datasets have training data available, so evaluation is in-distribution. *Out-of-Distribution:* Models are trained only on Calc., NQ, and MBPP (in-distribution, marked as ‡), then evaluated on all datasets (out-of-distribution for AIME24, AIME25, AMC, HumanEval, WebQ, LQA2, and GPQA). Zero-shot Router uses direct prompting without examples, while Few-shot Router uses prompting with examples. The best results are highlighted in **bold**.

Method	Math Reasoning			Code		Arith.	Common.		Logic	Sci.	Avg.
	AIME24	AIME25	AMC	Human.	MBPP <sup>‡</sup>	Calc. <sup>‡</sup>	NQ <sup>‡</sup>	WebQ	LQA2	GPQA	
<i>Closed-Source Models</i>											
Gemini2.5-Pro	92.0	86.7	62.5	81.5	83.7	64.7	59.2	63.5	78.9	84.0	75.6
GPT-5	93.3	94.6	97.5	93.4	98.4	82.9	59.3	61.5	83.8	85.7	85.0
GPT-4.1	46.7	33.3	82.5	92.1	57.7	62.0	54.5	61.5	78.2	62.1	63.0
GPT-4o	13.3	6.7	45.8	85.4	82.6	58.1	59.4	63.0	72.9	44.4	53.1
<i>Training-free Baselines</i>											
ZS Router	13.3	6.7	32.5	53.0	64.2	55.7	29.2	39.2	45.3	24.6	36.4
FS Router	23.3	13.3	40.0	68.9	64.7	47.2	27.3	35.8	40.8	25.9	38.7
Random Router	6.7	3.3	15.0	37.8	52.6	40.2	25.3	32.1	49.2	30.6	29.3
<i>In-Distribution Performance</i>											
ReAct (Yao et al., 2023)	13.3	13.3	37.5	64.6	71.4	48.5	31.3	38.4	45.3	35.3	39.9
Reflexion (Shinn et al., 2023)	13.3	13.3	40.0	61.0	73.4	52.1	34.6	36.8	46.7	33.9	40.5
RouterDC	40.0	23.3	62.5	80.5	77.7	74.9	41.2	47.6	47.2	39.1	53.4
MLPRouter	26.7	10.0	45.0	76.2	68.7	48.2	32.1	40.4	41.2	34.8	42.3
BertRouter	30.0	13.3	45.0	75.4	72.1	77.1	38.9	50.4	47.1	36.6	48.6
<b>ATLAS (cluster)</b>	<b>43.3</b>	<b>40.0</b>	<b>82.5</b>	<b>91.5</b>	<b>83.6</b>	<b>83.3</b>	<b>43.8</b>	<b>53.6</b>	<b>66.8</b>	<b>46.4</b>	<b>63.5</b>
<i>Out-of-Distribution Performance</i>											
ReAct	6.7	3.3	32.5	66.5	72.5	52.4	32.8	36.0	40.2	27.2	37.0
Reflexion	6.7	6.7	37.5	64.0	72.9	50.7	34.9	37.2	43.9	28.8	38.3
RouterDC	13.3	3.3	47.5	79.2	78.7	70.8	40.1	50.8	50.4	28.6	46.3
MLPRouter	13.3	3.3	32.5	75.0	67.7	54.6	37.3	43.7	38.9	26.8	39.3
BertRouter	6.7	6.7	40.0	78.7	79.0	67.0	38.9	51.4	40.3	27.7	43.6
<b>ATLAS (cluster)</b>	13.3	3.3	47.5	<b>91.5</b>	<b>83.6</b>	<b>83.3</b>	43.8	51.4	45.6	29.0	49.2
<b>ATLAS (RL)</b>	<b>43.3</b>	<b>33.3</b>	<b>67.5</b>	85.4	81.8	81.6	<b>44.1</b>	<b>52.2</b>	<b>62.7</b>	<b>42.0</b>	<b>59.4</b>

reasoning: ATLAS achieves 40.0% on AIME25 and 82.5% on AMC (+16.7% and +20.0% over RouterDC). Notably, ATLAS(cluster) exceeds GPT-4o (53.1%) and approaches GPT-4.1 (63.0%), demonstrating that strategic model–tool orchestration enables a reasoning engine of smaller-scale models to rival larger proprietary systems.

This performance stems from exploiting rich empirical priors through semantic embedding. By mapping queries into structured clusters and caching historical performance patterns, the framework achieves near-optimal task-configuration alignment. In contrast, supervised routers like BertRouter and MLPRouter struggle with non-linear decision boundaries in heterogeneous model-tool spaces. Their classification-based selection fails to capture nuanced synergies from domain-specific pairings, resulting in suboptimal routing.

#### 4.2.2 Generalization Scenarios

When facing out-of-distribution (OOD) challenges, ATLAS(cluster) suffers significant degradation (e.g., dropping from 40.0% to 3.3% on AIME25) as well as other baselines, whereas ATLAS(RL) maintains an average accuracy of 59.4% with 10.2% higher than ATLAS(cluster) (49.2%) and 13.1% higher

than RouterDC (46.3%). The gap is most striking in mathematical reasoning: on AIME24 and AIME25, ATLAS(RL) sustains 43.3% and 33.3% accuracy, respectively, while the clustering method achieves only 13.3% and 3.3% (a 10× difference). This indicates that the RL path learns transferable collaborative decision principles rather than task-specific mappings.

ATLAS(RL) autonomously explores effective trajectories through multi-faceted reward signals, learning generalizable patterns of model-tool synergies: when to invoke symbolic tools for verification or route to reasoning-specialized models rather than memorizing task-specific mappings. This enables robust transfer, maintaining competitive performance on unfamiliar tasks like AIME24 (43.3%) and GPQA (42.0%), approaching or exceeding GPT-4o despite using only 7B and 8B models. These results confirm that RL-driven component provides essential generalization capability, effectively bridging established domain expertise and unseen reasoning challenges.

#### 4.3 Multi-modal Tool Orchestration

To evaluate ATLAS on multi-modal tasks, we benchmark it against single-tool baselines across

Table 2: **Performance evaluation under dynamic routing pool extensions.** † denotes results after integrating domain-specialized models (Llama-3.1-8B-UltraMedical, Qwen2.5-Math-7B-Instruct) and an Outcome Reward Model into the routing pool. ‡ marks in-domain benchmarks; all others are out-of-domain. Best results are in **bold**.

Method	Math Reasoning			Code		Arith.	Common.		Logic	Sci.	Avg.
	AIME24	AIME25	AMC	Human.	MBPP‡	Calc.‡	NQ‡	WebQ	LQA2	GPQA	
ZS Router	13.3	6.7	32.5	53.0	64.2	55.7	29.2	39.2	45.3	24.6	36.4
ZS Router†	20.0	13.3	37.5	52.4	63.1	55.0	28.7	38.9	45.9	25.7	38.0
FS Router	23.3	13.3	40.0	68.9	64.7	47.2	27.3	35.8	40.8	25.9	38.7
FS Router†	26.7	16.7	47.5	70.7	63.8	46.5	25.9	36.2	41.7	25.0	40.0
RandomRouter	6.7	3.3	15.0	37.8	52.6	40.2	25.3	32.1	49.2	30.6	29.3
RandomRouter†	3.3	3.3	17.5	35.4	52.0	41.3	22.7	31.5	49.8	30.1	28.7
BertRouter	26.7	16.7	42.5	76.8	72.6	62.7	35.4	49.8	52.5	33.3	46.9
BertRouter†	33.3	20.0	50.0	75.0	73.0	61.3	36.2	50.1	53.4	32.4	48.4
ATLAS (RL)	43.3	33.3	67.5	<b>85.4</b>	<b>81.8</b>	81.6	44.1	52.2	62.7	42.0	59.4
<b>ATLAS (RL)†</b>	<b>50.0</b>	<b>40.0</b>	<b>70.0</b>	84.2	<b>81.8</b>	<b>82.4</b>	<b>45.3</b>	<b>52.8</b>	<b>64.8</b>	<b>45.1</b>	<b>61.7</b>

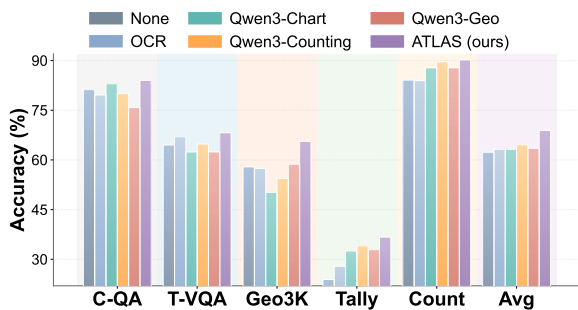


Figure 4: **Performance comparison of ATLAS against single-tool baselines across multi-modal benchmarks.** ‘None’ denotes direct reasoning without any tools. ATLAS achieves the highest accuracy.

five visual understanding and reasoning datasets, as shown in Figure 4. ATLAS achieves an average accuracy of 68.9% through dynamic tool invocation, outperforming the strongest single-tool baseline by 4.3%. Notably, ATLAS surpasses all individual tool in each task category. For example, exceeding the best single tool (Qwen3-Chart, 83.0%) on ChartQA and overcoming the performance limitations of single tools (e.g., Qwen3-Chart only achieves 50.2% on Geometry3K). This reveals that adaptive model-tool routing effectively integrates internal reasoning with external tool augmentation, thus establishing strong effectiveness on complex multi-modal tasks. Detailed results are provided in Appendix C.2.

#### 4.4 Generalization Toward Dynamic Model-Tool Synergy

A practical orchestration framework must accommodate an evolving ecosystem where new models and tools are continuously introduced. To evaluate this extensibility, we expand the routing pool

with three additional components: Llama-3.1-8B-UltraMedical (Zhang et al., 2024) for biomedical reasoning, Qwen2.5-Math-7B-Instruct (Yang et al., 2024b) for mathematical problem-solving, and an Outcome Reward Model for solution verification. Notably, our policy is trained exclusively on the original pool of 5 models and 4 tools; the newly added components are introduced only at inference time without any retraining. This extension substantially increases the combinatorial search space, posing a more challenging routing problem.

As shown in Table 2, ATLAS(RL) exhibits strong adaptability, improving from 59.4% to 61.7% (+2.3%) after pool extension. Gains are most pronounced on mathematical benchmarks: AIME24 (+6.7%) and AIME25 (+6.7%), confirming effective utilization of the newly added math-specialized model and verification tool. In contrast, baseline methods show limited or inconsistent responses: BertRouter gains only +1.5%, while RandomRouter degrades due to the expanded search space. This disparity arises because classifier-based routers learn fixed decision boundaries that become misaligned with new candidates, whereas ATLAS learns transferable routing principles through RL exploration, enabling seamless integration of new components without retraining.

#### 4.5 Discussion

**Evaluation on Reasoning Capacity Boundary.** Inspired by (Yue et al., 2025a), we implement the pass@k metric to measure the reasoning capacity boundary of ATLAS(RL), where pass@k equals 1 if at least one of k sampled outputs passes verification. As shown in Table 3, RL training yields

Table 3: **Reasoning capacity boundary analysis of ATLAS(RL)**. We report the pass@k metrics across diverse benchmarks to evaluate the exploration ( $k = 1$ ) and the potential reasoning upper bound ( $k = 16$ ).

	Math Reasoning			Code		Arith.	Common.		Logic	Sci.	Avg.
	AIME24	AIME25	AMC	Human.	MBPP <sup>‡</sup>	Calc. <sup>‡</sup>	NQ <sup>‡</sup>	WebQ	LQA2	GPQA	
<i>Pass@1 Results with/without ATLAS RL Training</i>											
w/o	13.3	6.7	32.5	53.0	64.2	55.7	29.2	39.2	45.3	24.6	36.4
w	<b>43.3</b>	<b>33.3</b>	<b>67.5</b>	<b>85.4</b>	<b>81.8</b>	<b>81.6</b>	<b>44.1</b>	<b>52.2</b>	<b>62.7</b>	<b>42.0</b>	<b>59.4</b>
△	+30.0	+26.6	+35.0	+32.4	+17.6	+25.9	+14.9	+13.0	+17.4	+17.4	+23.0
<i>Pass@16 Results with/without ATLAS RL Training</i>											
w/o	16.7	13.3	40.0	73.1	73.9	70.6	36.8	48.8	47.0	27.2	44.7
w	<b>50.0</b>	<b>36.7</b>	<b>75.0</b>	<b>89.6</b>	<b>84.5</b>	<b>83.3</b>	<b>46.9</b>	<b>54.9</b>	<b>64.4</b>	<b>45.8</b>	<b>63.1</b>
△	+33.3	+23.4	+35.0	+16.5	+10.6	+12.7	+10.1	+6.1	+17.4	+18.6	+18.4

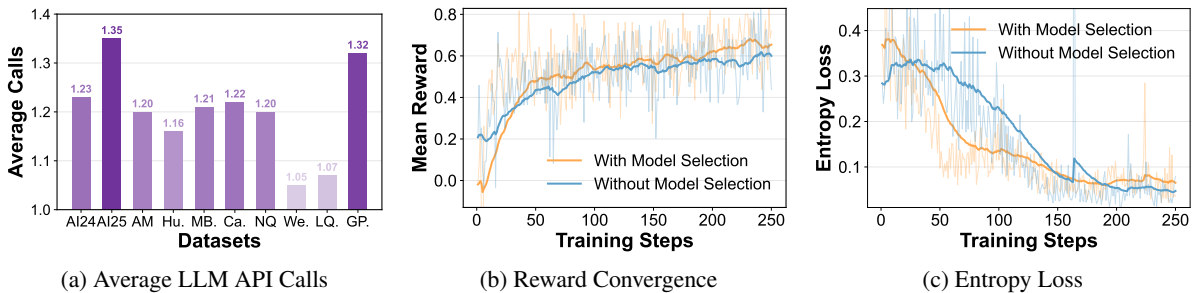


Figure 5: **Analysis of LLM API call count and ATLAS(RL) training dynamics.**

an absolute improvement of +23.0% in pass@1 accuracy (from 36.4% to 59.4%), demonstrating significantly optimized exploration efficiency. At pass@16, the upper bound reaches 63.1% (+3.7%), indicating that ATLAS(RL) already operates near its reasoning capacity ceiling, efficiently converging to optimal solutions without requiring extensive sampling. The trained model maintains substantial advantages across all tasks even at pass@16, with gains ranging from +6.1% to +35.0%, confirming that RL training effectively enhances agentic reasoning potential.

**Analysis of LLM API Call Count.** As illustrated in Figure 5a, ATLAS exhibits highly task-adaptive invocation patterns. For challenging reasoning-intensive tasks like AIME25 and GPQA, API calls increase significantly as the RL policy allocates higher computational budgets through multi-round routing and verification. Conversely, for straightforward retrieval tasks like WebQ and NQ, call counts remain minimal. This differentiated distribution confirms that ATLAS balances reasoning performance and inference cost, effectively suppressing redundant invocations where simpler models or fewer rounds suffice.

**Analysis of RL Training Convergence.** We validate the RL-driven routing policy’s stability

through reward and entropy evolution during training. Figure 5b shows that incorporating model selection reward ( $\mathcal{R}_{sel}$ ) yields faster convergence to a higher plateau compared to the baseline, guiding the agent toward higher-yield decision regions. Figure 5c demonstrates that the ATLAS configuration achieves a much sharper reduction in entropy compared to the ablation group, reaching a lower terminal value. This indicates that the router successfully transitions from stochastic exploration to a deterministic, high-confidence decision-making state, ensuring both the robustness and predictability of the routing process.

**Statistical Significance.** To verify that the observed gains are not due to chance, we conducted the Wilcoxon signed-rank test (Kotz and Johnson, 1992)—a nonparametric paired test that requires no normality assumption—by pairing ATLAS against the strongest baseline (RouterDC) across all benchmarks in Table 1. For in-distribution settings the test yields  $p = 9.7 \times 10^{-4}$  ( $p < \alpha=0.05$ ), decisively rejecting  $H_0$  and confirming that ATLAS’s superiority is statistically significant. Equivalent significance ( $p < 0.05$ ) is observed in the OOD setting. Variance analysis across three repeated runs further confirms stability: e.g., ATLAS(Cluster) achieves  $82.5 \pm 2.5$  on AMC and  $91.5 \pm 1.6$  on

HumanEval (in-distribution).

**More Discussion.** Due to space, we include more discussions in Appendix, including detailed multi-modal results (C.2), test-time scaling (C.3), analysis of model-tool preferences (C.4), ablation on reward design (C.7), and sensitivity analysis (C.8).

## 5 Conclusion

We present ATLAS, a generalizable framework for dynamic model-tool alignment through dual-path architecture: training-free cluster-based routing for domain-specific efficiency and RL-driven exploration for open-domain adaptability. ATLAS rivals or exceeds powerful closed-source models across diverse benchmarks, demonstrating a paradigm shift from model-centric scaling to ecosystem-centric orchestration. Experimental results show that strategic coordination of heterogeneous model-tool combinations unlocks superior reasoning while maintaining efficiency. As model and tool ecosystems continue to evolve, such orchestration reasoning systems will become essential for next-generation autonomous agents that address complex real-world challenges.

## Limitations

While ATLAS demonstrates strong performance across diverse benchmarks, several limitations warrant discussion. First, our current evaluation focuses primarily on text-based and visual reasoning tasks; extending to other modalities (e.g., audio, video) remains unexplored. Second, our framework assumes reliable API access to candidate models and tools-network latency or service unavailability in real-world deployments may impact performance. We plan to investigate more lightweight policy architectures and robust fallback mechanisms in future work.

## Acknowledgments

This work is supported by the National Natural Science Foundation of China (No. U2436210, No. 62322120).

## Ethical Considerations

All datasets, models, and tools utilized in this work are derived from publicly available resources with proper citations, involving no private or sensitive information. ATLAS consists of two components: a training-free cluster-based router and

an RL-trained policy model that learns to orchestrate existing LLMs and tools. While the policy model is trained to make routing decisions, the underlying candidate models and tools remain unmodified. Consequently, our framework inherits the potential biases, safety limitations, and ethical concerns present in these constituent components. ATLAS itself does not introduce new harmful capabilities beyond those already existing in the routing pool. We recommend that practitioners carefully evaluate all candidate models and tools for compliance with ethical guidelines, and apply appropriate safety measures when deploying ATLAS in real-world applications.

## References

- Manoj Acharya, Kushal Kafle, and Christopher Kanan. 2019. Tallyqa: Answering complex counting questions. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 8076–8084.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. Semantic parsing on Freebase from question-answer pairs. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1533–1544, Seattle, Washington, USA. Association for Computational Linguistics.
- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, et al. 2024. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*.
- Cong Chen, Kaixiang Ji, Hao Zhong, Muzhi Zhu, Anzhou Li, Guo Gan, Ziyuan Huang, Cheng Zou, Jiajia Liu, Jingdong Chen, et al. 2025a. Guishepherd: Reliable process reward and verification for long-sequence gui tasks. *arXiv preprint arXiv:2509.23738*.
- Guoxin Chen, Zhong Zhang, Xin Cong, Fangda Guo, Yesai Wu, Yankai Lin, Wenzheng Feng, and Yasheng Wang. 2025b. Learning evolving tools for large language models. In *The Thirteenth International Conference on Learning Representations*.
- Lingjiao Chen, Matei Zaharia, and James Zou. 2023. Frugalgpt: How to use large language models while reducing cost and improving performance. *Transactions on Machine Learning Research*.
- Mark Chen. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.

- Shuhao Chen, Weisen Jiang, Baijiong Lin, James Kwok, and Yu Zhang. 2024. Routerdc: Query-based router by dual contrastive learning for assembling large language models. *Advances in Neural Information Processing Systems*, 37:66305–66328.
- Dujian Ding, Ankur Mallick, Chi Wang, Robert Sim, Subhabrata Mukherjee, Victor Rühle, Laks V. S. Lakshmanan, and Ahmed Hassan Awadallah. 2024. Hybrid LLM: Cost-efficient and quality-aware query routing. In *The Twelfth International Conference on Learning Representations*.
- Guanting Dong, Yifei Chen, Xiaoxi Li, Jiajie Jin, Hongjin Qian, Yutao Zhu, Hangyu Mao, Guorui Zhou, Zhicheng Dou, and Ji-Rong Wen. 2025. Tool-star: Empowering llm-brained multi-tool reasoner via reinforcement learning. *arXiv preprint arXiv:2505.16410*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407.
- Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, Yujia Qin, Baoquan Zhong, Chengquan Jiang, Jinxin Chi, and Wanjun Zhong. 2025a. Retool: Reinforcement learning for strategic tool use in llms. *arXiv preprint arXiv:2504.11536*.
- Tao Feng, Yanzhen Shen, and Jiaxuan You. 2025b. Graphrouter: A graph-based router for LLM selections. In *The Thirteenth International Conference on Learning Representations*.
- Guo Gan, Yuxuan Ding, Cong Chen, Yuwei Ren, Yin Huang, and Hong Zhou. 2026. Android coach: Improve online agentic training efficiency with single state multiple actions. *arXiv preprint arXiv:2604.07277*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. In *The Eleventh International Conference on Learning Representations*.
- Qitian Jason Hu, Jacob Bieker, Xiuyu Li, Nan Jiang, Benjamin Keigwin, Gaurav Ranganath, Kurt Keutzer, and Shriyash Kaustubh Upadhyay. 2024. Routerbench: A benchmark for multi-LLM routing system. In *Agentic Markets Workshop at ICML 2024*.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. 2024. Qwen2.5-coder technical report. *arXiv preprint arXiv:2409.12186*.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan O Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025a. Search-r1: Training LLMs to reason and leverage search engines with reinforcement learning. In *Second Conference on Language Modeling*.
- Ruihan Jin, Pengpeng Shao, Zhengqi Wen, Jinyang Wu, Mingkuan Feng, Shuai Zhang, and Jianhua Tao. 2025b. Radialrouter: Structured representation for efficient and robust large language models routing. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 14587–14600, Suzhou, China. Association for Computational Linguistics.
- Yoonsik Kim, Moonbin Yim, and Ka Yeon Song. 2024. Tablevqa-bench: A visual question answering benchmark on multiple table domains. *arXiv preprint arXiv:2404.19205*.
- Yilun Kong, Jingqing Ruan, Yihong Chen, Bin Zhang, Tianpeng Bao, Shi Shiwei, Xiaoru Hu, Hangyu Mao, Ziyue Li, Xingyu Zeng, et al. 2024. Tptu-v2: Boosting task planning and tool usage of large language model-based agents in real-world industry systems. In *Proceedings of the 2024 conference on empirical methods in natural language processing: industry track*, pages 371–385.
- Samuel Kotz and Norman L. Johnson, editors. 1992. *Breakthroughs in Statistics: Methodology and Distribution*. Springer New York, New York, NY.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. 2023. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*.
- Yangning Li, Weizhi Zhang, Yuyao Yang, Wei-Chieh Huang, Yaozu Wu, Junyu Luo, Yuanchen Bei, Henry Peng Zou, Xiao Luo, Yusheng Zhao, Chunkit Chan, Yankai Chen, Zhongfen Deng, Yinghui Li, Haitao Zheng, Dongyuan Li, Renhe Jiang, Ming Zhang, Yangqiu Song, and Philip S. Yu. 2025. A survey of RAG-reasoning systems in large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 12120–12145, Suzhou, China. Association for Computational Linguistics.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*.

- Hanmeng Liu, Jian Liu, Leyang Cui, Zhiyang Teng, Nan Duan, Ming Zhou, and Yue Zhang. 2023. Logiqa 2.0—an improved dataset for logical reasoning in natural language understanding. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:2947–2962.
- Keming Lu, Hongyi Yuan, Runji Lin, Junyang Lin, Zheng Yuan, Chang Zhou, and Jingren Zhou. 2024. Routing to the expert: Efficient reward-guided ensemble of large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1964–1974.
- Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. 2021. Inter-GPS: Interpretable geometry problem solving with formal language and symbolic reasoning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6774–6786, Online. Association for Computational Linguistics.
- Zhengxi Lu, Yuxiang Chai, Yaxuan Guo, Xi Yin, Liang Liu, Hao Wang, Han Xiao, Shuai Ren, Pengxiang Zhao, Guangyi Liu, et al. 2026a. Ui-r1: Enhancing efficient action prediction of gui agents by reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 40, pages 17608–17616.
- Zhengxi Lu, Zhiyuan Yao, Jinyang Wu, Chengcheng Han, Qi Gu, Xunliang Cai, Weiming Lu, Jun Xiao, Yueting Zhuang, and Yongliang Shen. 2026b. Skill0: In-context agentic reinforcement learning for skill internalization. *arXiv preprint arXiv:2604.02268*.
- Zhengxi Lu, Jiabo Ye, Fei Tang, Yongliang Shen, Haiyang Xu, Ziwei Zheng, Weiming Lu, Ming Yan, Fei Huang, Jun Xiao, et al. 2025. Ui-s1: Advancing gui automation via semi-online reinforcement learning. *arXiv preprint arXiv:2509.11543*.
- Zexiong Ma, Chao Peng, Qunhong Zeng, Pengfei Gao, Yanzhen Zou, and Bing Xie. 2025. Tool-integrated reinforcement learning for repo deep search. *arXiv preprint arXiv:2508.03012*.
- MAA. 2024. [American invitational mathematics examination - aime](#). In *American Invitational Mathematics Examination - AIME 2024*.
- MAA. 2025. [American invitational mathematics examination - aime](#). In *American Invitational Mathematics Examination - AIME 2025*.
- Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the association for computational linguistics: ACL 2022*, pages 2263–2279.
- Kai Mei, Wujiang Xu, Shuhang Lin, and Yongfeng Zhang. 2025. Omnirouter: Budget and performance controllable multi-llm routing. *arXiv preprint arXiv:2502.20576*.
- Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E Gonzalez, M Waleed Kadous, and Ion Stoica. 2024. Routellm: Learning to route llms from preference data. In *The Thirteenth International Conference on Learning Representations*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel. 2023. Teaching clip to count to ten. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3170–3180.
- Zhihong Pan, Kai Zhang, Yuze Zhao, and Yupeng Han. 2025. Route to reason: Adaptive routing for llm and reasoning strategy selection. *arXiv preprint arXiv:2505.19435*.
- Shishir G Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. 2025. The berkeley function calling leaderboard (BFCL): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741.
- N Dinesh Reddy and Sudeep Pillai. 2025. Orion: A unified visual agent for multimodal perception, advanced visual reasoning and execution. *arXiv preprint arXiv:2511.14210*.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2024. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Yuhao Shen, Tianyu Liu, Junyi Shen, Jinyang Wu, Quan Kong, Li Huan, and Cong Wang. 2026a. Double: Breaking the acceleration limit via double retrieval speculative parallelism. *arXiv preprint arXiv:2601.05524*.

- Yuhao Shen, Junyi Shen, Quan Kong, Tianyu Liu, Yao Lu, and Cong Wang. 2026b. Specbranch: Speculative decoding via hybrid drafting and rollback-aware branch parallelism. In *The Fourteenth International Conference on Learning Representations*.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Hunyuan Vision Team, Pengyuan Lyu, Xingyu Wan, Gengluo Li, Shangpin Peng, Weinong Wang, Liang Wu, Huawen Shen, Yu Zhou, Canhui Tang, et al. 2025. Hunyuanocr technical report. *arXiv preprint arXiv:2511.19575*.
- Chenxu Wang, Hao Li, Yiqun Zhang, Linyao Chen, Jianhao Chen, Ping Jian, Peng Ye, Qiaosheng Zhang, and Shuyue Hu. 2025. Icl-router: In-context learned model representations for llm routing. *arXiv preprint arXiv:2510.09719*.
- Jinyang Wu, Mingkuan Feng, Shuai Zhang, Fangrui Lv, Ruihan Jin, Feihu Che, Zengqi Wen, and Jianhua Tao. 2025a. Boosting multimodal reasoning with automated structured thinking. *arXiv preprint arXiv:2502.02339*.
- Jinyang Wu, Chonghua Liao, Mingkuan Feng, Shuai Zhang, Zhengqi Wen, Haoran Luo, Ling Yang, Huazhe Xu, and Jianhua Tao. 2025b. Templaterl: Structured template-guided reinforcement learning for llm reasoning. *arXiv preprint arXiv:2505.15692*.
- Jinyang Wu, Shuo Yang, Changpeng Yang, Yuhao Shen, Shuai Zhang, Zhengqi Wen, and Jianhua Tao. 2026. Spark: Strategic policy-aware exploration via dynamic branching for long-horizon agentic learning. *arXiv preprint arXiv:2601.20209*.
- Jinyang Wu, Shuai Zhang, Feihu Che, Mingkuan Feng, Pengpeng Shao, and Jianhua Tao. 2025c. Pandora’s box or aladdin’s lamp: A comprehensive analysis revealing the role of RAG noise in large language models. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5019–5039, Vienna, Austria. Association for Computational Linguistics.
- Shirley Wu, Shiyu Zhao, Qian Huang, Kexin Huang, Michihiro Yasunaga, Kaidi Cao, Vassilis Ioannidis, Karthik Subbian, Jure Leskovec, and James Y Zou. 2024. Avatar: Optimizing llm agents for tool usage via contrastive reasoning. *Advances in Neural Information Processing Systems*, 37:25981–26010.
- Wenxun Wu, Yuanyang Li, Guhan Chen, Linyue Wang, and Hongyang Chen. 2025d. Tool-augmented policy optimization: Synergizing reasoning and adaptive tool use with reinforcement learning. *arXiv preprint arXiv:2510.07038*.
- Junlin Xie, Zhihong Chen, Ruifei Zhang, Xiang Wan, and Guanbin Li. 2024. Large multimodal agents: A survey. *arXiv preprint arXiv:2402.15116*.
- Fangzhi Xu, Hang Yan, Qiushi Sun, Jinyang Wu, Zixian Huang, Muye Huang, Jinyang Gong, Zichen Ding, Kanzhi Cheng, Yian Wang, et al. 2026. Odysseya: Benchmarking large language models for long-horizon, active and inductive interactions. *arXiv preprint arXiv:2602.05843*.
- Hang Yan, Xinyu Che, Fangzhi Xu, Qiushi Sun, Zichen Ding, Kanzhi Cheng, Jian Zhang, Tao Qin, Jun Liu, and Qika Lin. 2026. Tide: Trajectory-based diagnostic evaluation of test-time improvement in llm agents. *arXiv preprint arXiv:2602.02196*.
- Jianhao Yan, Yafu Li, Zican Hu, Zhi Wang, Ganqu Cui, Xiaoye Qu, Yu Cheng, and Yue Zhang. 2025. Learning to reason under off-policy guidance. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024a. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. 2024b. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2023. [React: Synergizing reasoning and acting in language models](#). In *The Eleventh International Conference on Learning Representations*.
- Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. 2025a. Does reinforcement learning really incentivize reasoning capacity in LLMs beyond the base model? In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Yanwei Yue, Guibin Zhang, Boyang Liu, Guancheng Wan, Kun Wang, Dawei Cheng, and Yiyan Qi. 2025b. MasRouter: Learning to route LLMs for multi-agent systems. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15549–15572, Vienna, Austria. Association for Computational Linguistics.
- Haozhen Zhang, Tao Feng, and Jiaxuan You. 2025a. Router-rl: Teaching llms multi-round routing and aggregation via reinforcement learning. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Kaiyan Zhang, Sihang Zeng, Ermo Hua, Ning Ding, Zhang-Ren Chen, Zhiyuan Ma, Haoxin Li, Ganqu Cui, Biqing Qi, Xuekai Zhu, Xingtai Lv, Hu Jinfang,

Zhiyuan Liu, and Bowen Zhou. 2024. Ultramedical: Building specialized generalists in biomedicine. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Yiqun Zhang, Hao Li, Chenxu Wang, Linyao Chen, Qiaosheng Zhang, Peng Ye, Shi Feng, Daling Wang, Zhen Wang, Xinrun Wang, Jia Xu, Lei Bai, Wanli Ouyang, and Shuyue Hu. 2025b. The avengers: A simple recipe for uniting smaller language models to challenge proprietary giants. *arXiv preprint arXiv:2505.19797*.

Richard Zhuang, Tianhao Wu, Zhaojin Wen, Andrew Li, Jiantao Jiao, and Kannan Ramchandran. 2025. EmbedLLM: Learning compact representations of large language models. In *The Thirteenth International Conference on Learning Representations*.

Yuxin Zuo, Kaiyan Zhang, Li Sheng, Shang Qu, Ganqu Cui, Xuekai Zhu, Haozhan Li, Yuchen Zhang, Xinwei Long, Ermo Hua, Biqing Qi, Youbang Sun, Zhiyuan Ma, Lifan Yuan, Ning Ding, and Bowen Zhou. 2025. TTRL: Test-time reinforcement learning. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.

C.10 Statistical Significance and Variance Analysis . . . . . 22

**D Case Study . . . . . 22**

**E Additional Discussion . . . . . 23**

E.1 Distinguishing ATLAS from Prior Routing and Tool Usage Methods . . . . . 23

E.2 Discussion on RL Reward Design . . . . . 24

E.3 When to Use Cluster-Based vs. RL-Based Routing . . . . . 24

E.4 Future Work . . . . . 24

**Contents**

**A Details of Methodology . . . . . 14**

A.1 Complete Algorithm Implementations . . . . . 14

A.2 Detailed Specification of Reward Signals . . . . . 14

**B Additional Experimental Details . . . . . 14**

B.1 Datasets . . . . . 14

B.2 Baselines . . . . . 17

B.3 Evaluation Details. . . . . 17

B.4 Implementation Details . . . . . 18

**C More Results and Analysis . . . . . 19**

C.1 Details Main Results . . . . . 19

C.2 Detailed Multimodal Results . . . . . 19

C.3 Test-Time Scaling Results . . . . . 19

C.4 Analysis of Model-Tool Alignment Preferences . . . . . 20

C.5 Scalability of RL-Based Routing with Tool Pool Size . . . . . 20

C.6 Evaluation on Realistic Tool-Calling Benchmark (BFCL) . . . . . 20

C.7 Ablation Study on Reward Components . . . . . 21

C.8 Sensitivity Analysis on Cluster Number . . . . . 21

C.9 Sensitivity of Cluster-Based Routing to Embedding Model . . . . . 22

## A Details of Methodology

### A.1 Complete Algorithm Implementations

We provide detailed implementations about training-free cluster-based routing (Algorithm 1), and RL-driven Multi-step Routing (Algorithm 2).

### A.2 Detailed Specification of Reward Signals

To bridge the gap between structured interaction and task-specific accuracy, ATLAS employs a composite reward function  $r_\phi = \mathcal{R}_{\text{fmt}} + \gamma\mathcal{R}_{\text{out}} + \xi\mathcal{R}_{\text{sel}}$ . This section provides the formal definitions and criteria for each reward component.

**Format Reward ( $\mathcal{R}_{\text{fmt}}$ )** The format reward ensures that the RL agent adheres to the predefined syntactic protocols, which is essential for stable parsing and environment interaction.  $\mathcal{R}_{\text{fmt}}$  is set to 0 if all the following conditions are satisfied, and  $-1$  otherwise:

- **Tag Integrity:** All XML-style tags (e.g., `<think>`, `<route>`, and `<answer>`) must be correctly opened and closed in a nested or sequential manner.
- **Invocation Syntax:** Tool calls within the search block must strictly follow the format `Model-Name@@Tool-Name:Input`. Furthermore, the specified model and tool names must exist within the active routing pool  $\mathcal{P}$ .
- **Mandatory Reasoning:** The trajectory must contain at least one complete `<think>...</think>` block to ensure internal deliberation before an action or answer.
- **Uniqueness of Response:** The trajectory must conclude with exactly one `<answer>...</answer>` block.
- **Execution Consistency:** To maintain the integrity of the multi-step interaction, the number of search calls initiated by the agent must strictly match the number of information blocks returned by the environment.

**Outcome Reward ( $\mathcal{R}_{\text{out}}$ )** The outcome reward serves as the primary signal for task success. It is a binary indicator evaluated upon the completion of the trajectory:

$$\mathcal{R}_{\text{out}} = \begin{cases} 1, & \text{if the answer } y_j \text{ is correct,} \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

**Model Selection Reward ( $\mathcal{R}_{\text{sel}}$ )** To encourage the agent to select the most efficient and capable expert for a given domain, we introduce an alignment-based penalty. The “optimal model” for each task is pre-determined as follows:

- For the **MBPP** dataset, the optimal model is defined as `Qwen2.5-Coder-7B-Instruct`.
- For the **Calculator** and **NQ** datasets, the optimal model is identified via an offline evaluation where GPT-4o judges the best-performing candidate from the pool for each specific query.

The reward is then formulated to penalize sub-optimal invocations:

$$\mathcal{R}_{\text{sel}} = \begin{cases} 0, & \text{if select the optimal model,} \\ -0.15, & \text{otherwise.} \end{cases} \quad (9)$$

## B Additional Experimental Details

### B.1 Datasets

The datasets utilized in this paper are summarized in Table 4. Below, we provide detailed descriptions of each benchmark to illustrate the diverse reasoning capabilities required by our framework.

#### Mathematical Reasoning.

- **AIME 2024 & AIME 2025 (MAA, 2024, 2025):** The American Invitational Mathematics Examination (AIME) is a prestigious 15-question, 3-hour test designed for high-performing high school students. We evaluate on the 2024 and 2025 editions, each containing 30 problems that demand advanced problem-solving skills, strategic thinking, and precise numerical computation.
- **AMC (Lightman et al., 2023):** The American Mathematics Competitions (AMC) consist of multiple-choice problems ranging from elementary to intermediate difficulty. Our evaluation set includes 40 problems that assess fundamental mathematical reasoning and computational proficiency.

#### Arithmetic Reasoning.

- **Calculator (Wu et al., 2025d):** A benchmark containing 1,000 complex arithmetic problems requiring precise numerical computation. These

---

**Algorithm 1** Training-free Cluster-based Routing

---

**Input:** Test query  $q_j$ , cluster centroids  $\{\mu_1, \dots, \mu_K\}$ , historical performance statistics  $\text{Stats}[k][m, t]$ , performance-cost trade-off parameter  $\alpha$ ;

**Output:** Optimal model-tool pair  $(m^*, t^*)$  and the generated response  $y_j$ ;

*// Step 1: Query Representation*

1:  $\mathbf{v}_j \leftarrow \text{Embed}(q_j)$  ▷ Project query into the latent embedding manifold

*// Step 2: Semantic Clustering*

2:  $k^* \leftarrow \arg \min_{k \in \{1, \dots, K\}} \text{dist}(\mathbf{v}_j, \mu_k)$  ▷ Find the nearest semantic cluster centroid

*// Step 3: Dynamic Selection*

3: **if**  $\text{Stats}[k^*]$  is not empty **then**

4:   **for** each candidate pair  $(m, t) \in \mathcal{M} \times \mathcal{T}$  **do**

5:      $\mathcal{U}_{k^*}^{(m,t)} \leftarrow (1 - \alpha) \cdot \text{Accuracy}(\text{Stats}[k^*]) - \alpha \cdot \text{Cost}(\text{Stats}[k^*])$

6:   **end for**

7:    $(m^*, t^*) \leftarrow \arg \max_{(m,t)} \mathcal{U}_{k^*}^{(m,t)}$  ▷ Select the optimal combination

8: **else**

9:    $(m^*, t^*) \leftarrow \text{FallbackStrategy}(q_j)$  ▷ Handle out-of-distribution queries

10: **end if**

11:  $y_j \leftarrow \text{Execute}(m^*, t^*, q_j)$  ▷ Invoke the selected model with the specific tool

12: **return**  $y_j$

---

problems test the model’s ability to recognize when external calculation tools are necessary and to correctly formulate and interpret computational results, evaluating the integration of reasoning and tool invocation.

Questions span diverse topics and require retrieving and synthesizing information from Wikipedia articles. This benchmark evaluates knowledge-intensive reasoning and information retrieval capabilities.

### Code Generation.

- **HumanEval** (Chen, 2021): This benchmark comprises 164 hand-crafted programming problems designed to evaluate code synthesis capabilities. Each problem includes a function signature, docstring, body, and unit tests. Solutions require understanding natural language specifications and generating functionally correct Python code.
- **MBPP** (Austin et al., 2021): The Mostly Basic Programming Problems (MBPP) dataset contains 974 crowd-sourced Python programming problems designed for entry-level programmers. Problems are described in natural language and require generating short Python functions, typically 1-10 lines of code. This benchmark tests basic programming constructs including loops, conditionals, and string manipulation.

### Commonsense Reasoning.

- **Natural Questions (NQ)** (Kwiatkowski et al., 2019): A question-answering dataset containing real user queries issued to Google Search.

- **Web Questions (WebQ)** (Berant et al., 2013): A dataset of 1,000 questions designed to test knowledge-based question answering. Questions are sourced from web search queries and require retrieving factual information from knowledge bases, evaluating the model’s ability to access and reason over external knowledge sources.

### Logical Reasoning.

- **LogiQA2** (Liu et al., 2023): An improved version of LogiQA containing 1,572 multiple-choice logical reasoning problems. Questions are sourced from standardized exams and require identifying logical relationships, drawing inferences, and evaluating argument structures. This benchmark tests formal reasoning capabilities including deductive, inductive, and abductive reasoning.

### Scientific Reasoning.

- **GPQA** (Rein et al., 2024): The Graduate-Level Google-Proof Q&A benchmark consists

---

**Algorithm 2** RL-driven Multi-step Routing

---

**Input:** Query  $q_j$ , policy  $\pi_\theta$ , reference  $\pi_{\text{ref}}$ , pool  $\mathcal{P}$ , parameters  $T_{\text{max}}, \theta, \beta, \gamma, \xi$ ;**Output:** Response  $y_j$  and trajectory  $\tau$ ;

```
// Step 1: Initialization
1:  $\tau \leftarrow \emptyset, C_0 \leftarrow \emptyset, s_0 \leftarrow \{q_j, C_0\}$ 
// Step 2: Multi-step Reasoning Loop
2: for  $t = 0$  to  $T_{\text{max}} - 1$  do
3:    $a_t \sim \pi_\theta(\cdot | s_t, \mathcal{P})$  ▷ Action  $a_t \in \{\text{think}, \text{route}(m, t_{\text{tool}})\}$ 
4:   if  $a_t = \text{think}$  then
5:      $o_t \leftarrow \pi_\theta.\text{Reasoning}(s_t)$  ▷ Internal reasoning
6:   else
7:      $o_t \leftarrow \text{Execute}(m, t_{\text{tool}}, s_t)$  ▷ Dynamic routing and tool invocation
8:   end if
9:    $C_{t+1} \leftarrow C_t \cup \{a_t, o_t\}, s_{t+1} \leftarrow \{q_j, C_{t+1}\}$ 
10:   $\tau \leftarrow \tau \cup \{(s_t, a_t, o_t)\}$ 
11:  if  $a_t$  contains Final Answer then break
12: end for
13:  $y_j \leftarrow \text{ParseAnswer}(\tau)$  ▷ Answer extraction
// Step 3: Policy Update (Training Mode)
14: if training_mode then
15:    $r_\phi(\tau) \leftarrow \mathcal{R}_{\text{fmt}} + \gamma \mathcal{R}_{\text{out}} + \xi \mathcal{R}_{\text{sel}}$ 
16:    $\mathcal{L}_\theta \leftarrow - \left[ r_\phi(\tau) \cdot \log \pi_\theta(\tau) - \beta \cdot \log \frac{\pi_\theta(\tau)}{\pi_{\text{ref}}(\tau)} \right]$ 
17:   Update  $\theta$  via PPO update rule:  $\nabla_\theta \mathcal{L}_\theta$ 
18: end if
19: return  $(y_j, \tau)$ 
```

---

Table 4: Detailed information on the datasets and test set sizes used in our experiments.

Category	Dataset	#Test Samples
Mathematical Reasoning	AIME 2024 (MAA, 2024)	30
	AIME 2025 (MAA, 2025)	30
	AMC (Lightman et al., 2023)	40
Code Generation	HumanEval (Chen, 2021)	164
	MBPP (Austin et al., 2021)	974
Arithmetic Reasoning	Calculator (Calc.) (Wu et al., 2025d)	1000
Commonsense Reasoning	Natural Questions (NQ) (Kwiatkowski et al., 2019)	1200
	Web Question (WebQ) (Berant et al., 2013)	1000
Logical Reasoning	LogiQA2 (Liu et al., 2023)	1572
Scientific Reasoning	GPQA (Rein et al., 2024)	448
Multi-modal Perception and Reasoning	ChartQA (Masry et al., 2022)	500
	Geometry3K (Lu et al., 2021)	601
	TallyQA (Acharya et al., 2019)	498
	CountBench (Paiss et al., 2023)	491
	TableVQA (Kim et al., 2024)	500

of 448 multiple-choice questions across biology, physics, and chemistry, written by domain

experts with PhD-level knowledge. Questions are designed to be difficult even for experts and

require deep domain understanding beyond simple fact retrieval.

### Multi-modal Reasoning.

- **ChartQA** (Masry et al., 2022): A visual question-answering benchmark containing questions about various chart types (bar charts, line graphs, pie charts). Questions require extracting quantitative information from visual representations and performing numerical reasoning, testing the integration of visual perception and mathematical computation.
- **Geometry3K** (Lu et al., 2021): A comprehensive geometry problem-solving dataset comprising multiple problems with diagram annotations. Problems involve diverse geometric concepts including angles, areas, perimeters, and spatial relationships. This benchmark evaluates visual-geometric reasoning and the ability to apply mathematical principles to diagrammatic representations.
- **TallyQA** (Acharya et al., 2019): A visual counting dataset containing complex counting questions across diverse real-world images. Questions range from simple object counting to complex scenarios requiring spatial reasoning and selective attention. This benchmark tests fine-grained visual perception and numerical reasoning capabilities.
- **CountBench** (Paiss et al., 2023): A specialized counting benchmark with questions designed to evaluate precise object enumeration in images. Unlike traditional counting tasks, CountBench emphasizes accuracy on challenging cases involving occlusions, similar objects, and cluttered scenes, requiring robust visual understanding.
- **TableVQA** (Kim et al., 2024): A visual question-answering benchmark containing questions about tables across multiple domains. Questions require understanding table structures, extracting relevant information, and performing reasoning over tabular data, evaluating the integration of visual perception and structured data comprehension.

These diverse benchmarks collectively assess the framework’s ability to dynamically select optimal model-tool combinations across varying task requirements, ranging from symbolic mathematical reasoning to multi-modal visual understanding.

## B.2 Baselines

In our experiments, we compare the proposed methods against six baseline approaches. Below, we provide detailed descriptions of each baselines.

- **Zero-shot (ZS) Router:** A baseline that directly prompts a base LLM to select the most suitable candidate model-tool combination from the available pool without prior examples.
- **Few-shot (FS) Router:** An extension of the zero-shot approach that incorporates several in-context examples to provide the base LLM with task-specific demonstrations and routing guidance.
- **Random Router:** A stochastic baseline that selects a candidate model-tool combination uniformly at random from the candidate pool for each query.
- **RouterDC** (Chen et al., 2024): A routing framework based on dual contrastive learning that maps queries and model-tool combinations into a shared embedding space. It utilizes sample-LLM and sample-sample contrastive losses to optimize query-model alignment and selects the optimal combination via cosine similarity.
- **MLPRouter** (Hu et al., 2024): A classification-based framework that trains an MLP for each model-tool combination. Each MLP predicts the success probability of its corresponding combination, and the one with the highest output is selected.
- **BertRouter** (Ong et al., 2024): A router utilizing a pre-trained mDeBERTaV3-base encoder (He et al., 2021) with an integrated classification head to predict the accuracy of model-tool pairings, following a selection logic similar to MLPRouter.

## B.3 Evaluation Details.

Our experiments employ two evaluation protocols: **In-Distribution (ID)**, where each dataset has its own training split, and **Out-of-Distribution (OOD)**, where models are trained exclusively on three datasets (Calculator, NQ, MBPP) and evaluated on all ten benchmarks, making AIME24, AIME25, AMC, HumanEval, WebQ, LogiQA2, and GPQA fully out-of-domain. For cluster-based

routing in OOD settings, semantic clusters and performance statistics are derived *solely* from the three training datasets; test queries from unseen domains are assigned to the nearest cluster based on semantic similarity, without accessing any OOD test set information. This design reflects realistic domain-specific scenarios but inevitably suffers from cluster misalignment on unfamiliar tasks (49.2% OOD vs. 63.5% ID, Table 1). In contrast, RL-based routing learns transferable patterns (when to invoke symbolic tools or defer to specialized models) that generalize beyond training distributions, achieving 59.4% OOD accuracy. Importantly, no test set information is leaked: all routing decisions rely purely on query embeddings and training domain statistics, ensuring evaluation integrity and demonstrating that gains stem from our dual-path architecture’s complementary strengths.

#### B.4 Implementation Details

##### Hyperparameters for Cluster-based Routing.

We set the number of cluster centers to 8 and employ the KMeans algorithm with the following hyperparameters: the cluster centers are initialized using the k-means++ method to accelerate convergence; the algorithm is allowed up to 1000 iterations per run; the number of initializations is set to automatic selection, the hyperparameter  $\alpha$  in Equation 4 is set to 0.5; and the Elkan variant of KMeans is used for computational efficiency.

**Hyperparameters for RL Training.** We train the policy model (Qwen2.5-3B-Instruct) using PPO with generalized advantage estimation (GAE). The training and validation batch sizes are both set to 24. The maximum prompt length is 4096 tokens, while the maximum response length is set to 3000 tokens. To control context growth, the maximum lengths for the observations are set to 2048 tokens each, and the maximum number of interaction turns is limited to 4. The actor is optimized with a learning rate of  $1 \times 10^{-6}$ , while the critic uses a learning rate of  $1 \times 10^{-5}$ . The PPO mini-batch size and micro-batch size for the actor are set to 12 and 6, respectively. The KL-divergence coefficient is fixed to 0.001. During rollout, we use a temperature of 1.0. For the reward weights in  $r_\phi = \mathcal{R}_{\text{fmt}} + \gamma \mathcal{R}_{\text{out}} + \xi \mathcal{R}_{\text{sel}}$ , we assign  $\gamma = \xi = 1$ . All experiments are conducted for 250 total training steps. We also provide the RL system prompt in Figure 6.

**Tool Details.** We provide the system prompt for three special multimodal tools in Figure 7-9. Re-

garding text-based tools, the *Code Interpreter* executes Python code, returns the execution results, indicates whether the execution was successful, and reports error locations and underlying causes in case of failures. The *Web Search* tool leverages the official Google Custom Search API to retrieve the three most relevant search result snippets. Search results are obtained by sending HTTP GET requests to the API (<https://www.googleapis.com/customsearch/v1>) with the required parameters, including the API key, search engine ID, query string, and the number of top results to return. The *Calculator* parses the model output in a function-call format to extract the mathematical expression, computation type, and precision requirements, and then computes and returns the result using appropriate functions from the *sympy* library. The *Process Reward Model (PRM)* runs five model outputs in parallel, evaluates the segmented outputs using a reward model, and selects the output with the highest average score as the final result. In this work, we use the off-the-shelf Qwen2.5-Math-PRM-7B<sup>1</sup>.

**Baseline Details.** The original baselines perform routing among multiple models. When evaluating these baselines, we replace the models in the code with model–tool combinations, thereby enabling the baselines to route over both models and tools simultaneously. For example, when reproducing EmbedLLM (Zhuang et al., 2025), we substitute the model names in the training data with the model–tool combination names, and replace the model performance with the empirically measured performance of the model–tool combinations. Apart from these modifications, all training and evaluation procedures strictly follow the official open-source implementations, ensuring that the reported results faithfully reflect the true performance of the baseline methods. When evaluating closed-source models, we use exactly the same evaluation code and prompt templates (including the use of CoT reasoning and fixed answer formats) as those used for other baselines and our proposed method, ensuring a strictly fair comparison and convincing final results.

**Computing Details.** All experiments are conducted on eight NVIDIA A100-80GB GPUs.

<sup>1</sup>Qwen/Qwen2.5-Math-PRM-7B

Table 5: **Extended performance comparison across diverse tasks and domains.** *In-Distribution:* All datasets have training data available, so evaluation is in-distribution. *Out-of-Distribution:* Models are trained only on Calc., NQ, and MBPP (in-distribution, marked as ‡), then evaluated on all datasets (out-of-distribution for AIME24, AIME25, AMC, HumanEval, WebQ, LQA2, and GPQA). Zero-shot Router uses direct prompting without examples, while Few-shot Router uses prompting with examples. The best results are highlighted in **bold**.

Method	Math Reasoning			Code		Arith.	Common.		Logic	Sci.	Avg.
	AIME24	AIME25	AMC	Human.	MBPP <sup>‡</sup>	Calc. <sup>‡</sup>	NQ <sup>‡</sup>	WebQ	LQA2	GPQA	
<i>Closed-Source Models</i>											
Gemini2.5-Pro	92.0	86.7	62.5	81.5	83.7	64.7	59.2	63.5	78.9	84.0	75.6
Gemini2.5-Flash	88.0	78.0	72.5	80.5	82.6	58.9	54.9	61.3	74.6	78.3	73.0
GPT-5	93.3	94.6	97.5	93.4	98.4	82.9	59.3	61.5	83.8	85.7	85.0
GPT-4.1	46.7	33.3	82.5	92.1	57.7	62.0	54.5	61.5	78.2	62.1	63.0
GPT-4o	13.3	6.7	45.8	85.4	82.6	58.1	59.4	63.0	72.9	44.4	53.1
<i>Training-free Baselines</i>											
ZS Router	13.3	6.7	32.5	53.0	64.2	55.7	29.2	39.2	45.3	24.6	36.4
FS Router	23.3	13.3	40.0	68.9	64.7	47.2	27.3	35.8	40.8	25.9	38.7
Random Router	6.7	3.3	15.0	37.8	52.6	40.2	25.3	32.1	49.2	30.6	29.3
<i>In-Distribution Performance</i>											
RouterDC	40.0	23.3	62.5	80.5	77.7	74.9	41.2	47.6	47.2	39.1	53.4
GraphRouter	30.0	16.7	50.0	78.7	75.0	72.3	37.5	49.6	45.8	37.1	49.3
EmbedLLM	23.3	13.3	45.0	75.6	72.0	76.7	36.8	48.3	51.4	35.0	47.7
MLPRouter	26.7	10.0	45.0	76.2	68.7	48.2	32.1	40.4	41.2	34.8	42.3
BertRouter	30.0	13.3	45.0	75.4	72.1	77.1	38.9	50.4	47.1	36.6	48.6
<b>ATLAS (cluster)</b>	<b>43.3</b>	<b>40.0</b>	<b>82.5</b>	<b>91.5</b>	<b>83.6</b>	<b>83.3</b>	<b>43.8</b>	<b>53.6</b>	<b>66.8</b>	<b>46.4</b>	<b>63.5</b>
<i>Out-of-Distribution Performance</i>											
RouterDC	13.3	3.3	47.5	79.2	78.7	70.8	40.1	50.8	50.4	28.6	46.3
GraphRouter	16.7	3.3	42.5	76.2	73.4	71.2	36.5	49.3	47.2	27.7	44.4
EmbedLLM	13.3	3.3	45.0	79.9	73.0	79.1	41.4	50.2	51.5	31.7	46.8
MLPRouter	13.3	3.3	32.5	75.0	67.7	54.6	37.3	43.7	38.9	26.8	39.3
BertRouter	6.7	6.7	40.0	78.7	79.0	67.0	38.9	51.4	40.3	27.7	43.6
<b>ATLAS (cluster)</b>	<b>13.3</b>	<b>3.3</b>	<b>47.5</b>	<b>91.5</b>	<b>83.6</b>	<b>83.3</b>	<b>43.8</b>	<b>51.4</b>	<b>45.6</b>	<b>29.0</b>	<b>49.2</b>
<b>ATLAS (RL)</b>	<b>43.3</b>	<b>33.3</b>	<b>67.5</b>	85.4	81.8	81.6	<b>44.1</b>	<b>52.2</b>	<b>62.7</b>	<b>42.0</b>	<b>59.4</b>

## C More Results and Analysis

### C.1 Details Main Results

We provide extended comparisons in Table 5, incorporating additional closed-source models (e.g., Gemini-2.5-Flash) and routing baselines including GraphRouter (Feng et al., 2025b) and EmbedLLM (Zhuang et al., 2025).

### C.2 Detailed Multimodal Results

Visual perception, comprehension, and reasoning are crucial capabilities for autonomous agents (Xie et al., 2024; Wu et al., 2025a; Reddy and Pillai, 2025). We conduct multimodal extension experiments as described in Section 4.3, with detailed orchestration results provided in Table 6. Our evaluation spans diverse visual reasoning tasks: chart understanding (ChartQA (Masry et al., 2022), TableVQA (Kim et al., 2024)), math reasoning (Geometry3K (Lu et al., 2021)), and object enumeration (TallyQA (Acharya et al., 2019), CountBench (Paiss et al., 2023)).

To ensure fair comparison, all configurations, in-

cluding the baseline (direct inference with Qwen3-8B-VL), single-tool baselines, and ATLAS, share the same foundational model (Qwen3-8B-VL) and identical evaluation protocols. The only modification across settings is the inclusion of tool-invocation instructions in the system prompt, which guide the model on when and how to invoke specific tools (e.g., chart parsing, object counting, geometric reasoning). Crucially, the core reasoning capacity and model parameters remain unchanged, ensuring that observed performance gains stem from adaptive tool orchestration rather than model-level differences or prompt engineering artifacts. The results demonstrate that ATLAS consistently outperforms single-tool baselines across all categories, validating the effectiveness of dynamic tool orchestration in multimodal scenarios. We plan to explore more backbones in future work.

### C.3 Test-Time Scaling Results

We analyze the scalability of our approach by increasing the self-consistency (SC) sample count on several representative benchmarks. As illus-

Table 6: **Performance comparison of ATLAS against single-tool baselines across multi-modal benchmarks.** The framework dynamically routes queries among multi-modal tools using Qwen3-VL-8B-Instruct as the backbone. ‘None’ represents direct reasoning without any tools. The best results are highlighted in **bold**.

Tool	Chart Understanding		Math Reasoning	Object Enumeration		Avg.
	ChartQA	TableVQA	Geometry3K	TallyQA	CountBench	
None (Direct Reasoning)	81.2	64.5	57.9	23.9	84.1	62.3
+OCR	79.6	67.0	57.4	27.9	83.9	63.2
+Qwen3-Chart	83.0	62.4	50.2	32.5	87.8	63.2
+Qwen3-Counting	80.0	64.8	54.4	34.1	89.6	64.6
+Qwen3-Geo	75.8	62.4	58.7	32.9	87.8	63.5
<b>ATLAS (ours)</b>	<b>84.0</b>	<b>68.2</b>	<b>65.6</b>	<b>36.7</b>	<b>90.2</b>	<b>68.9</b>
$\Delta$ vs. None	+2.8	+3.7	+7.7	+12.8	+6.1	+6.6

Table 7: Performance scaling with Self-Consistency (SC) across different sample sizes.

Dataset	Pass@1	SC@4	SC@8	SC@16
AIME24	43.3	63.3	66.7	70.0
AIME25	40.0	43.3	46.7	50.0
AMC	82.5	92.5	95.0	97.5
Calc.	83.3	83.5	84.7	86.9
GPQA	46.4	53.6	57.1	59.4

trated in Table 7, performance across almost all datasets shows a positive correlation with the number of samples ( $k$ ). For example, on the AIME24 benchmark, SC@16 yields a significant improvement from 43.3% to 70.0%. Similar findings are also observed on other tasks, such as commonsense reasoning and scientific reasoning. These results demonstrate that the ensemble of model-tool combinations provides a more robust candidate pool for majority voting.

#### C.4 Analysis of Model-Tool Alignment Preferences

Table 8 illustrates the strategic alignment between specific models and tools across diverse benchmarks. In deterministic domains such as coding (HumanEval) and advanced mathematics (AIME), ATLAS exhibits a clear convergence, selecting specialized pairings like Qwen2.5-Coder-7B with Python or DeepSeek-R1 with PRM in nearly 100% of cases. This high degree of consistency confirms the framework’s ability to internalize the performance advantages of domain-specific modules.

In contrast, knowledge-intensive tasks (NQ, MedQA) trigger a transition toward retrieval-augmented configurations, primarily utilizing Llama-3.1-8B with Web-Search. For more complex, broad-spectrum benchmarks like LQA2, the selection distribution becomes significantly more

granular, with the dominant combination of ATLAS(RL) accounting for only 56.4%. This shift demonstrates that ATLAS avoids rigid heuristics, instead employing a flexible orchestration strategy that adapts to the specific nuances and difficulty of each query.

#### C.5 Scalability of RL-Based Routing with Tool Pool Size

As shown in Table 9, we evaluate how the RL routing policy scales as the tool pool grows from 2 to 8 tools:

- **2 Tools:** Web-search, Python-code.
- **4 Tools:** Above + Calculator, PRM.
- **8 Tools:** Above + Reasoning Chain Refiner, Tool-Use Backtracer, Code Debugger, Segment Refiner.

As tool pool size grows from 4 to 8, accuracy improves by +1.5% on average while latency increases by only  $\sim 2$  s, indicating **sub-linear latency growth**. This demonstrates that the RL policy effectively identifies relevant tools even as the action space expands, owing to the transferable routing principles learned during training.

#### C.6 Evaluation on Realistic Tool-Calling Benchmark (BFCL)

We evaluate ATLAS(RL) on the Berkeley Function Calling Leaderboard (BFCL) (Patil et al., 2025), a benchmark specifically designed for realistic, production-oriented tool-invocation scenarios. As shown in Table 10, we report single-round results across four sub-categories.

ATLAS(RL) achieves an average score of 86.0%, surpassing all listed closed-source models and validating the effectiveness of our approach in complex, real-world tool-calling scenarios.

Table 8: **Distribution of dominant model-tool combinations across diverse benchmarks.** Dominant combination indicates the most frequently selected model-tool pair by our framework for each specific dataset.

Dataset	Dominant Combination	ATLAS(Cluster)	Dominant Combination	ATLAS(RL)
AIME24	DeepSeek.-7B@PRM	100.0%	DeepSeek.-7B@PRM	100.0%
AIME25	DeepSeek.-7B@PRM	100.0%	DeepSeek.-7B@PRM	100.0%
AMC	DeepSeek.-7B@PRM	100.0%	DeepSeek.-7B@PRM	91.7%
Human.	Coder-7B@Python	100.0%	Coder-7B@Python	100.0%
MBPP	Coder-7B@Python	100.0%	Coder-7B@Python	100.0%
Calc.	Qwen2.5-7B@Calc.	100.0%	Qwen2.5-7B@Calc.	95.8%
NQ	Llama3.1-8B@Search	92.8%	Llama3.1-8B@Search	99.0%
WebQ	Llama3.1-8B@Search	98.8%	Llama3.1-8B@Search	100.0%
LQA2	InternLM3-8B@Search	99.7%	InternLM3-8B@Search	56.4%
GPQA	DeepSeek.-7B@Python	80.4%	DeepSeek.-7B@Python	95.5%

Table 9: **Performance and latency under different tool pool sizes.**

Number of Tools	AIME25	AMC	MBPP	Calc.	LQA2	GPQA	Avg.
<i>Accuracy (%)</i>							
2	23.3	60.0	72.3	56.1	60.3	40.2	52.0
4	33.3	67.5	81.8	81.6	62.7	42.0	61.5
8	<b>36.7</b>	<b>70.0</b>	<b>81.8</b>	<b>83.3</b>	<b>62.9</b>	<b>42.9</b>	<b>62.9</b>
<i>Time Latency per Query (s)</i>							
2	67.96	34.16	3.91	0.34	1.28	75.96	30.6
4	123.15	45.20	4.02	0.26	1.35	84.68	43.1
8	132.62	44.84	4.85	0.42	1.65	86.95	45.2

### C.7 Ablation Study on Reward Components

To investigate individual reward contributions, we train the RL policy without  $\mathcal{R}_{\text{sel}}$  or  $\mathcal{R}_{\text{fmt}}$  while keeping other signals intact. This addresses concerns about potential circularity from GPT-4o judgments in  $\mathcal{R}_{\text{sel}}$  and validates the necessity of format enforcement.

As shown in Table 11, removing  $\mathcal{R}_{\text{sel}}$  causes a modest 3.1% degradation (59.4%  $\rightarrow$  56.3%), with notable drops on mathematical reasoning (AIME24/25:  $-6.6\%$  each). However, the policy still substantially outperforms all baselines, including RouterDC (46.3%) and EmbedLLM (46.8%). The retained performance (56.3% vs. 36.4% for zero-shot) confirms that ATLAS learns effective routing independently through  $\mathcal{R}_{\text{fmt}}$  and  $\mathcal{R}_{\text{out}}$ , without requiring external model judgments. This validates  $\mathcal{R}_{\text{sel}}$  as an efficiency-oriented auxiliary signal rather than a necessary component.

In contrast, removing  $\mathcal{R}_{\text{fmt}}$  leads to a more substantial 6.1% degradation (59.4%  $\rightarrow$  53.3%), with significant drops on mathematical reasoning (AIME24:  $-10.0\%$ , AMC:  $-12.5\%$ ) and code generation (HumanEval:  $-7.4\%$ ). This reveals that format enforcement is critical for maintaining structured interaction patterns—proper tool syntax and reasoning-action sequencing—which form the foundation for multi-step orchestration. Without

$\mathcal{R}_{\text{fmt}}$ , the policy produces malformed tool calls that propagate failures throughout reasoning trajectories. These results validate our design:  $\mathcal{R}_{\text{fmt}}$  and  $\mathcal{R}_{\text{out}}$  constitute essential signals, while  $\mathcal{R}_{\text{sel}}$  provides optional efficiency guidance.

### C.8 Sensitivity Analysis on Cluster Number

To evaluate the robustness of cluster-based routing to the choice of cluster granularity, we conduct sensitivity analysis by varying the number of clusters  $K \in \{4, 8, 16\}$  while keeping all other hyperparameters fixed. As shown in Table 12, the optimal performance is achieved at  $K = 8$  with 63.5% average accuracy, representing an 11.6% improvement over  $K = 4$  (51.9%) and a modest 0.7% gain over  $K = 16$  (62.8%). The substantial performance drop at  $K = 4$  suggests that overly coarse clustering fails to capture fine-grained task distinctions, leading to suboptimal model-tool alignments, particularly evident on code generation (HumanEval: 43.3% vs. 91.5%) where diverse programming patterns require more specialized routing. Conversely, increasing to  $K = 16$  yields diminishing returns, as excessively fine-grained clusters may suffer from data sparsity within each partition, resulting in less reliable performance statistics. These results demonstrate that moderate cluster granularity ( $K = 8$ ) strikes an effective balance

Table 10: **BFCL single-round tool invocation results.**

Method	Simple	Multiple	Parallel	Multiple Parallel	Avg.
o3	74.3	89.0	86.5	78.0	81.9
GPT-5.2	72.9	88.0	89.0	77.5	81.9
Grok-4	67.0	93.5	89.0	87.5	84.3
Gemini-2.5-Flash	74.3	92.0	94.0	79.5	85.0
<b>ATLAS(RL)</b>	<b>76.0</b>	<b>93.5</b>	<b>91.0</b>	<b>83.5</b>	<b>86.0</b>

Table 11: **Ablation study on reward components.** We evaluate the impact of removing  $\mathcal{R}_{\text{sel}}$  (model selection reward) and  $\mathcal{R}_{\text{fmt}}$  (format reward) on out-of-distribution performance. Models are trained on Calc., NQ, and MBPP ( $\ddagger$ ), then evaluated on all datasets. The best results are highlighted in **bold**.

Method	Math Reasoning			Code		Arith.	Common.		Logic	Sci.	Avg.
	AIME24	AIME25	AMC	Human.	MBPP $\ddagger$	Calc. $\ddagger$	NQ $\ddagger$	WebQ	LQA2	GPQA	
<i>Training-free Baselines</i>											
ZS Router	13.3	6.7	32.5	53.0	64.2	55.7	29.2	39.2	45.3	24.6	36.4
FS Router	23.3	13.3	40.0	68.9	64.7	47.2	27.3	35.8	40.8	25.9	38.7
Random Router	6.7	3.3	15.0	37.8	52.6	40.2	25.3	32.1	49.2	30.6	29.3
<i>Training-based Baselines</i>											
RouterDC	13.3	3.3	47.5	79.2	78.7	70.8	40.1	50.8	50.4	28.6	46.3
GraphRouter	16.7	3.3	42.5	76.2	73.4	71.2	36.5	49.3	47.2	27.7	44.4
EmbedLLM	13.3	3.3	45.0	79.9	73.0	79.1	41.4	50.2	51.5	31.7	46.8
MLPRouter	13.3	3.3	32.5	75.0	67.7	54.6	37.3	43.7	38.9	26.8	39.3
BertRouter	6.7	6.7	40.0	78.7	79.0	67.0	38.9	51.4	40.3	27.7	43.6
<b>ATLAS (RL)</b>	<b>43.3</b>	<b>33.3</b>	<b>67.5</b>	<b>85.4</b>	<b>81.8</b>	<b>81.6</b>	<b>44.1</b>	<b>52.2</b>	62.7	<b>42.0</b>	<b>59.4</b>
w/o $\mathcal{R}_{\text{sel}}$	36.7	26.7	65.0	82.3	80.6	79.1	41.3	48.3	<b>62.9</b>	40.6	56.3
$\Delta$	-6.6	-6.6	-2.5	-3.1	-1.2	-2.5	-2.8	-3.9	+0.2	-1.4	-3.1
<b>ATLAS (RL)</b>	<b>43.3</b>	<b>33.3</b>	<b>67.5</b>	<b>85.4</b>	<b>81.8</b>	<b>81.6</b>	<b>44.1</b>	<b>52.2</b>	<b>62.7</b>	<b>42.0</b>	<b>59.4</b>
w/o $\mathcal{R}_{\text{fmt}}$	33.3	26.7	55.0	78.0	75.4	78.3	41.6	48.0	58.2	38.4	53.3
$\Delta$	-10.0	-6.6	-12.5	-7.4	-6.4	-3.3	-2.5	-4.2	-4.5	-3.6	-6.1

between semantic specificity and statistical robustness, though the framework remains reasonably stable across a range of  $K$  values (62.8%–63.5% for  $K \in \{8, 16\}$ ), indicating limited sensitivity to this hyperparameter in practical deployments.

### C.9 Sensitivity of Cluster-Based Routing to Embedding Model

As shown in Table 13, we evaluate five embedding models spanning 300M–7B parameters to assess whether the cluster-based routing is sensitive to encoder choice.

Performance variance across models is minimal (63.1%–63.9%), demonstrating that the cluster-based routing is **fundamentally robust to embedding model choice**. This stability arises because semantic clustering depends primarily on the quality of the semantic embedding space rather than encoder-specific characteristics: as long as the encoder captures task-level semantics, the downstream cluster-to-performance mapping remains effective.

### C.10 Statistical Significance and Variance Analysis

To verify result stability, we repeat each experiment three times with different random seeds and report mean  $\pm$  standard deviation on three representative benchmarks in Table 14 and 15.

Beyond variance, we apply the Wilcoxon signed-rank test (Kotz and Johnson, 1992) pairing ATLAS against ROUTERDC across all 10 benchmarks in Table 1. The in-domain test yields  $p = 9.7 \times 10^{-4}$  ( $p < 0.05$ ), rejecting the null hypothesis  $H_0$  of no significant difference. The OOD setting yields equivalent significance ( $p < 0.05$ ), validating both the robustness and generalization capability of our framework.

### D Case Study

We provide some representative examples in Figures 10–14 to illustrate how ATLAS dynamically orchestrates model-tool combinations across diverse reasoning tasks.

Table 12: **Sensitivity analysis on cluster number  $K$ .** Performance across different cluster granularities in cluster-based routing. All datasets have training data available (in-distribution setting).

#Cluster	Math Reasoning			Code		Arith.	Common.		Logic	Sci.	Avg.
	AIME24	AIME25	AMC	Human.	MBPP <sup>‡</sup>	Calc. <sup>‡</sup>	NQ <sup>‡</sup>	WebQ	LQA2	GPQA	
4	36.7	30.0	75.0	43.3	71.5	79.1	28.8	48.5	66.8	39.6	51.9
8	<b>43.3</b>	<b>40.0</b>	<b>82.5</b>	<b>91.5</b>	<b>83.6</b>	<b>83.3</b>	<b>43.8</b>	<b>53.6</b>	<b>66.8</b>	<b>46.4</b>	<b>63.5</b>
16	40.0	<b>40.0</b>	<b>82.5</b>	90.9	82.9	82.3	44.1	53.4	66.7	45.3	62.8

Table 13: In-distribution performance across embedding models (cluster-based routing).

Embedding Model	AIME25	AMC	MBPP	LQA2	GPQA	Avg.
embedding-gemma-300m	40.0	82.5	83.8	66.2	46.2	63.7
jina-embeddings-v3	40.0	80.0	83.7	66.2	45.5	63.1
nomi-embed-text-v1.5	40.0	82.5	83.7	66.1	43.1	63.1
gte-Qwen2-1.5B-instruct	40.0	82.5	83.8	66.3	46.9	63.9
gte-Qwen2-7B-instruct (ours)	40.0	82.5	83.6	66.8	46.4	63.9

Table 14: **In-distribution performance comparison** (mean  $\pm$  std, three runs).

Method	AIME2025	AMC	HumanEval
RouterDC	22.2 $\pm$ 1.9	62.4 $\pm$ 4.3	80.5 $\pm$ 1.6
MLPRouter	8.9 $\pm$ 1.9	45.0 $\pm$ 2.5	76.0 $\pm$ 1.9
BertRouter	13.3 $\pm$ 3.4	44.2 $\pm$ 1.4	75.6 $\pm$ 1.2
<b>ATLAS(Cluster)</b>	<b>38.9 <math>\pm</math> 1.9</b>	<b>82.5 <math>\pm</math> 2.5</b>	<b>91.5 <math>\pm</math> 1.6</b>

Table 15: **Out-of-distribution performance comparison** (mean  $\pm$  std, three runs).

Method	AIME2025	AMC	HumanEval
RouterDC	2.2 $\pm$ 1.9	46.7 $\pm$ 1.4	79.1 $\pm$ 0.9
MLPRouter	3.3 $\pm$ 3.4	32.5 $\pm$ 2.5	75.1 $\pm$ 1.9
BertRouter	7.8 $\pm$ 1.9	39.2 $\pm$ 1.4	78.7 $\pm$ 1.6
ATLAS(Cluster)	3.3 $\pm$ 3.4	47.5 $\pm$ 2.5	<b>91.5 <math>\pm</math> 0.9</b>
<b>ATLAS(RL)</b>	<b>33.3 <math>\pm</math> 3.4</b>	<b>67.5 <math>\pm</math> 2.5</b>	85.4 $\pm$ 0.6

**Adaptive Multi-turn Reasoning.** Figure 10 demonstrates ATLAS’s capacity for self-correction through iterative exploration. When addressing a logical reasoning problem, the policy initially selects Qwen2.5-7B with web search to verify its hypothesis (option C), but upon receiving contradictory feedback, it re-evaluates the alternatives and routes to InternLM3-8B for a second verification. This multi-turn deliberation ultimately leads to the correct answer (option D), showcasing the framework’s ability to recover from suboptimal initial decisions through adaptive re-routing.

**Task-Aware Model-Tool Alignment and Selection.** Figures 11–14 highlight how ATLAS aligns model-tool pairs with task-specific requirements. For arithmetic computation (Figure 11), the policy directly invokes the calculator tool without unnecessary reasoning steps. For factual retrieval (Fig-

ure 12), it routes to Llama-3.1-8B with web search, recognizing the need for external knowledge. Code generation tasks (Figure 13) are delegated to the specialized Qwen2.5-Coder model with Python execution. For challenging mathematical problems (Figure 14), ATLAS combines DeepSeek-7B with PRM for rigorous verification. These examples collectively demonstrate that ATLAS has internalized meaningful associations between task categories and optimal model-tool configurations, rather than relying on rigid heuristics.

## E Additional Discussion

### E.1 Distinguishing ATLAS from Prior Routing and Tool Usage Methods

While ATLAS employs established techniques such as semantic clustering for query representation and PPO for policy optimization, its contribution extends beyond the individual components to address a fundamental gap in existing literature: the joint optimization of heterogeneous model-tool combinations. Prior routing methods (Chen et al., 2024; Ong et al., 2024; Lu et al., 2024) focus exclusively on model selection, treating LLMs as isolated execution units without considering external tool augmentation. Conversely, tool usage frameworks (Feng et al., 2025a; Wu et al., 2025d) rely on fixed invocation logic that cannot dynamically adapt to different model capabilities. ATLAS unifies these two paradigms by explicitly modeling the Cartesian product space  $\mathcal{S} = \mathcal{M} \times \mathcal{T}$  and learning task-aware alignments within this joint space.

The technical novelty of ATLAS manifests in three key aspects. **First**, the dual-path architec-

ture strategically combines training-free cluster-based routing for exploiting domain-specific priors with RL-driven exploration for generalizing to unfamiliar tasks, achieving complementary strengths across distribution shifts (Table 1). **Second**, our composite reward structure ( $\mathcal{R}_{\text{fmt}} + \gamma\mathcal{R}_{\text{out}} + \xi\mathcal{R}_{\text{sel}}$ ) decouples execution correctness from routing efficiency through the  $\mathcal{R}_{\text{sel}}$  signal, enabling the policy to internalize transferable expertise distribution rather than memorizing task-specific mappings, which is evidenced by robust generalization to expanded model-tool pools without retraining (Section 4.4). **Third**, our controlled experiments ensure that all configurations share identical backbone models and evaluation protocols, with only routing mechanisms differing (Section B.4), thereby isolating the contribution of orchestration strategies from confounding factors such as model capacity or prompt engineering. The consistent performance gains across 15 benchmarks, including out-of-distribution settings (+13.1% over baselines) and multi-modal tasks (+4.3%), demonstrate that ATLAS captures generalizable principles for adaptive model-tool coordination.

## E.2 Discussion on RL Reward Design

Our composite reward function  $r_\phi = \mathcal{R}_{\text{fmt}} + \gamma\mathcal{R}_{\text{out}} + \xi\mathcal{R}_{\text{sel}}$  balances structured execution, task correctness, and routing efficiency. Regarding potential concerns about  $\mathcal{R}_{\text{out}}$  that require ground-truth labels, we note that test-time reinforcement learning remains effective in label-scarce scenarios through alternative supervision signals: majority voting across sampled trajectories has proven effective as pseudo-labeling (Zuo et al., 2025). Future extensions of ATLAS could integrate such self-verification mechanisms to further reduce reliance on explicit supervision.

Regarding  $\mathcal{R}_{\text{sel}}$ , which penalizes suboptimal model selections based on offline evaluation (Equation 9), a potential concern is whether this introduces evaluator bias or test-time information leakage. However,  $\mathcal{R}_{\text{sel}}$  encodes domain priors from offline profiling, such as “code tasks benefit from specialized models” or “retrieval tasks favor search-augmented models”, which practitioners naturally possess and use to initialize routing systems. Critically, it does not leak test-time information but rather provides consistent training targets to guide efficiency-aware exploration. The low weight  $\|\xi\mathcal{R}_{\text{sel}}\| = 0.15$  (vs.  $\|\gamma\mathcal{R}_{\text{out}}\| = 1.0$  for  $\mathcal{R}_{\text{out}}$ ) ensures routing efficiency serves as an aux-

iliary signal without overriding correctness. Our ablation (Figure 5b-5c) shows that  $\mathcal{R}_{\text{sel}}$  accelerates convergence and reduces entropy, while generalization to expanded model pools without prior annotations (Table 2) shows that the policy learns transferable routing principles that aligns task characteristics with model capabilities, rather than memorizing specific mappings.

## E.3 When to Use Cluster-Based vs. RL-Based Routing

The choice between cluster-based and RL-based routing depends on data availability and generalization requirements. When domain-specific training data is accessible, such as historical query-answer pairs in enterprise QA systems, cluster-based routing offers a simple and efficient solution. It achieves strong in-domain performance (63.5% average accuracy, Table 1) with zero training cost by leveraging semantic clustering and historical statistics, making it ideal for rapid deployment in well-defined domains. Conversely, when the reasoning engine must handle diverse, unfamiliar tasks where domain priors are unavailable, such as general-purpose assistants facing unpredictable queries, RL-based routing provides superior generalization. It learns transferable patterns of when to invoke tools or defer to specialized models, maintaining robust OOD performance (59.4% vs. 49.2% for cluster-based) at the cost of upfront training. In practice, practitioners can adopt a hybrid strategy: using cluster-based routing as the default for efficiency while reserving RL-based routing for critical queries or new domains, thereby balancing simplicity with adaptability.

## E.4 Future Work

Several promising directions remain for future exploration. First, the current web-search tool does not account for retrieval noise; integrating noise-aware retrieval strategies (Li et al., 2025; Wu et al., 2025c) could further improve the robustness of information-augmented reasoning. Second, the RL routing policy could benefit from structured template guidance (Yan et al., 2025; Wu et al., 2025b) to impose stronger inductive biases over multi-step trajectories, potentially accelerating convergence and improving sample efficiency. Third, evaluating ATLAS in long-horizon, interactive settings such as ODYSSEYARENA (Xu et al., 2026; Wu et al., 2026; Lu et al., 2026b; Yan et al., 2026) would provide a more rigorous assessment of its agentic capabil-

ities beyond static benchmarks. On the efficiency side, speculative parallelism techniques (Shen et al., 2026a,b) offer a natural avenue to reduce the multi-step inference latency incurred by the RL-driven routing path. Finally, extending the framework to richer perceptual abilities—such as unified visual agents that tightly couple perception with reasoning (Reddy and Pillai, 2025; Lu et al., 2025, 2026a; Chen et al., 2025a; Gan et al., 2026)—remains an important direction toward building truly general-purpose orchestration systems.

## System Prompt For RL Experiments

Answer the given question.

Every time you receive new information, you must first conduct reasoning inside `<think> ... </think>` to plan how to solve the problem and analyze which model and tool to use.

After reasoning, you must call a model and a tool for assistance by writing your request inside `<search> ... </search>`.

### TOOL SELECTION GUIDELINES

1. calculator:

- Use for simple or complex numerical calculations. But it is the only tool that need strictly function-call style interactions.

*[Omit some tool descriptions]...*

### AVAILABLE MODELS

1. llama-3.1-8b-instruct

- Best for common-sense questions and factual question answering that need a wide range of information.
- Choose this model when you need to solve common-sense reasoning and factual Q&A problems

*[Omit some model descriptions]...*

### !!! STRICT FORMAT RULES for <search> !!!

You MUST use the following format EXACTLY:

`<search> Model-Name@@Tool-Name: Your-Input </search>`

Example: `<search> llama-3.1-8b-instruct@@web-search: What is the capital of France? </search>`

...

### REASONING REQUIREMENTS BEFORE TOOL CALL

Before EACH `<search>` call, you MUST think inside `<think> ... </think>` about:

- + Which category of the problem is the question belongs to, like mathematical, factual QA, code generation, complex calculation, etc.
- + Why external information / calculation / verification / code generation is needed.
- + Which tool is best suited for answering it, based on the tools' abilities.

...

### TOOL RESPONSE HANDLING

When a tool is called, the response will be returned between:

`<information> ... </information>`

You may:

- + Call different tools multiple times.
- + Use different models for different tool calls.
- + Re-call the same tool if it remains the most effective.

### FINAL ANSWER FORMAT

If no further tool assistance is required, provide the final answer directly using:

`<answer> ... </answer>`

You must start with `<think>`.

Question: {question}

Figure 6: System prompt for ATLAS RL Experiments.

## System Prompt For Qwen3-Chart Tool

### [Step 1]

Extract the useful data for solving the question in the image. Do not give me the final answer. Focus on the number. Make sure the number is right.

### [Step 2]

The image is one of a bar chart, line chart, or pie chart. It comes with instructions, predicted information, a question, and an answer. Based on the image and its predicted information, the question must be solved.

Your response should contain four fields and four fields only:

- "instruction explanation": an explanation the process of following instructions according to the chart type
- "explanation": an explanation of how you arrived at the answer
- "answer": the answer based on the chart and question

#### Instruction for whole chart:

1. If the question is about values corresponding to specific positions (ex: lower, last, top), then you must match the information with the chart image's positions for reasoning.
2. If the question demands interpretation based on the magnitude of values, reasoning should be based on the information's values.
3. Originally, this task requires solving based solely on the image, meaning all positions should be interpreted based on the image itself.
4. In most cases, the presence of x-axis values or y-axis values enables the determination of the chart's values.
5. Note that, you can utilize the predicted information. The predicted columns and rows are very likely to correspond to the actual columns and rows of the chart, and this can help you determine where the rows and columns exist in the chart image.

#### Instruction for bar chart:

1. Firstly, bars of the same color represent the same column. Therefore, distinguishing colors and identifying corresponding columns is crucial (usually displayed around the main chart in the form of a legend).
2. Next, determine the location of rows. For vertical bar charts, rows are typically annotated at the bottom of the main chart, while for horizontal bar charts, they are annotated on the left or right side of the main chart.
3. Then, combine the colors of the nearest bars with annotated rows to determine which row and column the bars correspond to in the information.
4. Afterwards, locate the values corresponding to each row and column. If values are annotated on the bars, refer to them. Otherwise, compare the sizes of the bars to find the values.
5. For vertical bar charts, the y-axis value at the end of the bar corresponds to its value. Similarly, for horizontal bar charts, the x-axis value at the end of the bar corresponds to its value.

#### Instruction for line chart:

1. In the case of a line chart, the bottom x-axis will primarily represent the rows, and each colored line will represent a column.
2. The legend, which indicates which column corresponds to the color of the line, is usually located within the main chart. If the legend is absent or located separately, the text annotated with the color corresponding to the line will likely indicate the column (if colors are not present, the text annotated at the left or right end of the line is likely to correspond to the column).
3. The point of the line passing through the same x coordinates as each x-axis is the value itself (meaning the x-axis corresponds to the row, the color of the line corresponds to the column, and that point is the value).
4. If there is an annotation near a line point, it is highly likely that this value represents the value of the point.
5. If there is no annotation near a line point, you can determine the value of the point by referring to the y-axis value corresponding to the y coordinate of the point.
6. In a line chart, it is crucial to understand the flow of the line. Lines can show trends of decreasing, increasing, or remaining constant, and when multiple lines intersect, it is important to identify which line corresponds to which column based on their colors.

#### Instruction for pie chart:

1. In a pie chart, it is very important to determine which color corresponds to which row.
2. Each section has one color, and the row it corresponds to is likely indicated by text either inside the section or close to it (if not nearby, it can be identified through the legend or connected to the corresponding text by lines or markers).
3. In the case of a pie chart, the values are usually annotated on each section of the pie chart.

Now you have to answer the question based on the first image, predicted information and question.

Here is predicted information: {response\_of\_step\_1}

Here is a question: {question}

Figure 7: System prompt for Qwen3-Chart Tool.

## System Prompt For Qwen3-Counting Tool

### [Step 1]

List all the objects the question asked. Give me their exact position. Note that some objects may be occluded; do not overlook these partially obscured objects.

### [Step 2]

Answer the question based on the position information of the objects and the image. The position information may be incomplete; you need to identify objects that might have been missed.

#### # Example1:

There is a cup in the upper left corner of the image, and another cup at the bottom of the image.

**Question:** How many cups in the image?

**Response:** We first check whether any existing object position information is missing. Based on the image, there is indeed one cup in the upper-left corner and another at the bottom, so the answer is `\boxed{{2}}`.

#### # Example2:

There is a cup in the upper left corner of the image, and another cup at the bottom of the image.

**Question:** How many cups in the image?

**Response:** We first check whether any existing object location information is missing. Based on the image, in addition to the cups in the upper-left corner and at the bottom, there is also one in the upper-right corner, so the answer is `\boxed{{3}}`. The number can not be zero.

Question: {question}

Position information: {response\_of\_step\_1}

Figure 8: System prompt for Qwen3-Counting Tool.

## System Prompt For Qwen3-Geo Tool

### [Step 1]

Question: {question}

Extract ALL visual and text information in this structured format:

#### 1. POINTS:

- List each labeled point and its type (vertex, intersection, center, etc.)
- Format: "Label: [type] at approximate coordinates (x,y relative)"

#### 2. LINES & SEGMENTS: For each line/segment:

- Endpoints, Type (solid, dashed, dotted)
- Special markings (arrow for ray, tick marks for equal length)
- Labels if present

#### 3. ANGLES: For each marked angle:

- Vertex point
- Arms (two points forming the angle)
- Measurement if given
- Right angle/square symbol if present

#### 4. CIRCLES/ARCS:

- Center point.
- Radius/diameter if labeled
- Arc endpoints if partial

#### 5. ANNOTATIONS:

- Length measurements (with units)
- Angle measurements (with degree symbol)
- Ratio markings (single/double tick marks)
- Textual labels not attached to points

#### 6. RELATIONSHIPS:

- Parallel lines (mark with // )
- Perpendicular lines (mark with  $\perp$ )
- Tangency points
- Collinearity (points on same line)

Output this as a structured list with clear headings.

### [Step 2]

Image analyze: {response\_of\_step\_1}

Question: {question}

Begin step-by-step reasoning with self-explanation:

Step 1: [State current goal]

- Why: [Justification for this step]
- How: [Method/Theorem applied]
- Result: [Intermediate conclusion]
- Check: [Verify this step is valid]

Step 2: [Next goal based on Step 1 result]

- Why: ...

Continue until final result.

Apply geometric theorems explicitly: "If [condition], then by [theorem name], we know [conclusion]"

### [Step 3]

Entire solution: {response\_of\_step\_2}

Question: {question}

Check the entire solution for consistency:

#### 1. DIMENSIONAL ANALYSIS:

- Do units match throughout?
- Are angle measures in consistent units?
- Do ratios make sense?

#### 2. EXTREME CASE TESTING:

- Does solution hold for special configurations?
- What if a measurement approaches zero/infinity?
- Boundary condition check

#### 3. CROSS-METHOD VERIFICATION:

- Can result be obtained via alternative method?
- Quick mental estimation vs. calculated result
- Symmetry considerations

#### 4. PLAUSIBILITY CHECK: Is magnitude reasonable?

### [Step 4]

Entire solution: {response\_of\_step\_2}

Question: {question}

Test results: {response\_of\_step\_3}

If the reasoning is right, summary the final answer from the Reasoning steps. If it is wrong, correct it based on the Tests results. Remember to put your final answer on the last line using the format  $\boxed{\$ANSWER}$ .

Figure 9: System prompt for Qwen3-Counting Tool.

## Example 1 on LQA2

**Question** The radiation absorbed by someone during an ordinary commercial airline flight is no more dangerous than that received during an ordinary dental X-ray. Since a dental X-ray does negligible harm to a person, we can conclude that the radiation absorbed by members of commercial airline flight crews will also do them negligible harm. A flaw in the argument is its failure to consider that

- A. there may be many forms of dangerous radiation other than X-rays and the kinds of radiation absorbed by members of commercial airline flight crews
- B. receiving a dental X-ray may mitigate other health risks, whereas flying does not
- C. flying at high altitude involves risks in addition to exposure to minor radiation
- D. the longer and the more often one is exposed to radiation, the more radiation one absorbs and the more seriously one is harmed

**Ground truth** D

---

### Model output

<think>

The argument presented assumes that since dental X-rays do negligible harm, then exposure to radiation during commercial airline flights will also do negligible harm. However, this conclusion does not account for the cumulative effect of radiation exposure over time or the potential for additional health risks associated with flying. The argument fails to consider the possibility that prolonged exposure to low levels of radiation could accumulate to a harmful level, or that flying might expose individuals to other health risks in addition to radiation. This makes option C the most relevant flaw in the argument.

</think>

<search> qwen2.5-7b-instruct@@web-search: A flaw in the argument is its failure to consider that flying at high altitude involves risks in addition to exposure to minor radiation. Is this correct? </search>

<information> This statement is incorrect. </information>

<think>

The provided search result indicates that the statement "flying at high altitude involves risks in addition to exposure to minor radiation" is incorrect according to the web search. Therefore, option C cannot be the correct answer. Given the options, the argument's flaw must lie elsewhere. Let's reconsider the options:\n- Option A suggests considering other forms of dangerous radiation, but the argument does not claim that dental X-rays are the only form of radiation.\n- Option B mentions mitigating other health risks, but the argument does not discuss any other health risks.\n- Option D addresses the cumulative effect of radiation exposure and the potential for harm with increased exposure, which aligns with the argument's failure to consider long-term effects. Option D is the most appropriate choice because it directly addresses the cumulative nature of radiation exposure and the potential for harm with increased exposure, which the argument overlooks.

</think>

<search> internlm3-8b-instruct@@web-search: A flaw in the argument is its failure to consider that the longer and the more often one is exposed to radiation, the more radiation one absorbs and the more seriously one is harmed. Is this correct?

</search>

<information> Yes, this is correct. </information>

<answer> Answer: D </answer>

Figure 10: Example 1 on the LQA2 dataset.

### Example 2 on calculator

**Question** What is (-20 divided by -73) minus 12.256132902293373?

**Ground truth** -11.98

---

#### Model output

<think>

The given question is a numerical calculation that requires a calculator to perform the division and subtraction. There is no need for reasoning with the provided models or tools since a calculator is the most appropriate tool for this type of calculation.

</think>

<search> qwen2.5-7b-instruct@@calculator: The full original question is: What is (-20 divided by -73) minus 12.256132902293373? </search>

<information>Result: -11.982160299553646254366867651697</information>

<answer>Answer: -11.982160299553646</answer>

Figure 11: Example 2 on the Calculator dataset.

### Example 3 on WebQ

**Question** Which country was Justin Bieber born in?

**Ground truth** Canada

---

#### Model output

<think>

The birthplace of an individual can often be found through basic research or through public records. Since this is a factual question that does not require numerical calculation, coding, or external search, it is best suited for a model that can answer common-sense questions and factual questions, such as llama-3.1-8b-instruct.

</think>

<search> llama-3.1-8b-instruct@@web-search:Who is Justin Bieber and what country was he born in? </search>

<information>Justin Bieber is a Canadian singer. He was born in Canada.</information>

<answer>Answer: Canada</answer>

Figure 12: Example 3 on the WebQ dataset.

## Example 4 on MBPP

**Question** Write a function to find whether an array is subset of another array.

**Test cases** `assert is_subset([11, 1, 13, 21, 3, 7], 6, [11, 3, 7, 1], 4) == True`  
`assert is_subset([1, 2, 3, 4, 5, 6], 6, [1, 2, 4], 3) == True`  
`assert is_subset([10, 5, 2, 23, 19], 5, [19, 5, 3], 3) == False`

### Model output

**<think>**

The given task requires writing a function in Python. There is no need for any external information, numerical calculation, or verification. A Python code generation is exactly what is needed here.

**</think>**

**<search>** Coder-7B-Instruct@@python-code: Write a function to find whether an array is subset of another array.

**<information>**

```
def is_subset(arr1, m, arr2, n):
    set_arr1 = set(arr1)
    for i in range(n):
        if arr2[i] not in set_arr1:
            return False
    return True
```

**</information>**

**<answer>**

```
def is_subset(arr1, m, arr2, n):
    set_arr1 = set(arr1)
    for i in range(n):
        if arr2[i] not in set_arr1:
            return False
    return True
```

**</answer>**

Figure 13: Example 4 on the MBPP dataset.

## Example 5 on AIME25

**Question** Find the sum of all integer bases  $b > 9$  for which  $17_{(b)}$  is a divisor of  $97_{(b)}$ .

**Ground truth** 70

### Model output

<think>

To solve this problem, we need to convert the given numbers from base  $b$  to base 10 and then determine the conditions under which  $17_{(b)}$  is a divisor of  $97_{(b)}$ . The number  $17_{(b)}$  in base  $b$  converts to  $1 \cdot b + 7 = b + 7$  in base 10. Similarly,  $97_{(b)}$  converts to  $9 \cdot b + 7 = 9b + 7$  in base 10. For  $17_{(b)}$  to be a divisor of  $97_{(b)}$ , there must exist an integer  $k$  such that  $9b + 7 = k(b + 7)$ . Rearranging this equation, we get  $9b + 7 = kb + 7k$ , which simplifies to  $9b - kb = 7k - 7$  or  $b(9 - k) = 7(k - 1)$ . This implies that  $b$  must be a divisor of  $7(k - 1)$  and  $9 - k$  must be a positive integer. Since  $b > 9$ , we can test values of  $b$  to find which ones satisfy this condition.

</think>

<search> DeepSeek-R1-Distill-Qwen-7B@@prm: The original question is: Find the sum of all integer bases  $b > 9$  for which  $17_{(b)}$  is a divisor of  $97_{(b)}$ . </search>

<information>

To solve the problem, we need to find the sum of all integer bases  $(b > 9)$  for which  $(17_{(b)})$  is a divisor of  $(97_{(b)})$ .

First, we convert the numbers from base  $(b)$  to decimal:  $(17_{(b)} = 1 \cdot b + 7 = b + 7) - (97_{(b)} = 9 \cdot b + 7 = 9b + 7)$

We need  $(b + 7)$  to divide  $(9b + 7)$ . This can be written as:  $[\frac{9b + 7}{b + 7}]$

Performing polynomial long division, we find:  $[\frac{9b + 7}{b + 7} = 9 - \frac{56}{b + 7}]$

For this to be an integer,  $(\frac{56}{b + 7})$  must be an integer, meaning  $(b + 7)$  must be a divisor of 56.

The divisors of 56 are 1, 2, 4, 7, 8, 14, 28, and 56. Since  $(b > 9)$ ,  $(b + 7)$  must be greater than 16. The relevant divisors are 28 and 56.

Thus, we have:

If  $(b + 7 = 28)$ , then  $(b = 21)$ .

If  $(b + 7 = 56)$ , then  $(b = 49)$ .

Both 21 and 49 are valid bases. Summing these bases, we get:  $[21 + 49 = 70]$

Therefore, the sum of all such integer bases  $(b)$  is  $(\boxed{70})$ .

</information>

<answer>\boxed{70}</answer>

Figure 14: Example 5 on the AIME dataset.