# JoyAgents-R1: Joint Evolution Dynamics for Versatile Multi-LLM Agents with Reinforcement Learning

**Ai Han**[*†]**, Junxing Hu**[*]**, Pu Wei, Zhiqian Zhang, Yuhang Guo, Jiawei Lu, Zicheng Zhang**

JD.com

Beijing, China

## Abstract

Multi-agent reinforcement learning (MARL) has emerged as a prominent paradigm for increasingly complex tasks. However, joint evolution across heterogeneous agents remains challenging due to cooperative inefficiency and training instability. In this paper, we propose the joint evolution dynamics for MARL called JoyAgents-R1, which first applies Group Relative Policy Optimization (GRPO) to the joint training of heterogeneous multi-agents. By iteratively refining agents' large language models (LLMs) and memories, the method achieves holistic equilibrium with optimal decision-making and memory capabilities. Specifically, JoyAgents-R1 first implements node-wise Monte Carlo sampling on the behavior of each agent across entire reasoning trajectories to enhance GRPO sampling efficiency while maintaining policy diversity. Then, our marginal benefit-driven selection strategy identifies top-$K$ sampling groups with maximal reward fluctuations, enabling targeted agent model updates that improve training stability and maximize joint benefits through cost-effective parameter adjustments. Meanwhile, JoyAgents-R1 introduces an adaptive memory evolution mechanism that repurposes GRPO rewards as cost-free supervisory signals to eliminate repetitive reasoning and accelerate convergence. Experiments across general and domain-specific scenarios demonstrate that JoyAgents-R1 achieves performance comparable to that of larger LLMs while built on smaller open-source models.

## 1 Introduction

The rapid advancement of Large Language Models (LLMs) [1–5] has revolutionized agent-based systems, empowering agents with the ability to perform reasoning, planning, and natural language interaction across diverse domains [6]. Compared to single agents with specialized functionalities [7–10], multi-agent systems demonstrate superior flexibility and scalability in tackling complicated tasks such as Artificial Intelligence (AI) assistants [11] and emergency responder management [12], among which Multi-Agent Reinforcement Learning (MARL) is a mainstream subfield of the community [13].

Since Reinforcement Learning (RL) has demonstrated remarkable efficacy in aligning models with human preferences [14], LLM-based MARL methods have flourished and achieved certain results in sophisticated task decomposition [15, 16] and adaptive coordination [17, 18]. In MARL, the behavior of one agent may affect the rewards of other agents, which may cause environmental instability and lead to low system efficiency and performance [19]. While methods like Multi-Agent Proximal Policy Optimization (MAPPO) [20] have advanced MARL by adapting PPO [21] to multi-agent settings, their reliance on additional value functions introduces critical limitations. Moreover, the decoupling of policy and value updates in actor-critic architecture often leads to training instability, particularly when coordinating heterogeneous agents with misaligned reward structures [22], which poses severe challenges to the dynamics of multi-agent evolution.

---

[*]Equal contribution. [†]Corresponding author: <hanai5@jd.com>.

Recently, DeepSeek-R1 [23] introduced the Group Relative Policy Optimization (GRPO) [24], a novel RL framework that enhances decision-making in LLMs. Instead of relying on a critic model, GRPO uses population dynamics to generate multiple responses per input and select actions based on relative group advantages, significantly reducing computational overhead. Since GRPO has demonstrated significant performance gains in single or homogeneous agents [25–27] and vision-language models [28], we are inspired to explore its potential in general multi-agent joint evolution. However, directly applying GRPO to train general multi-agents which are typically heterogeneous encounters challenges. Since the joint action space grows exponentially with the number of agents and creates far more complex action spaces than those in single-agent scenarios [29, 30], tailored sampling and updating strategies are required to guide multi-agent evolution. In addition, heterogeneous multi-agents struggle to train and converge due to dynamic reasoning paths and diverse architectures.

To address the above issues, we propose JoyAgents-R1, a novel joint evolutionary framework that leverages GRPO to enable heterogeneous multi-agent systems to achieve the optimal equilibrium. To the best of our knowledge, this is the first work to apply GRPO in general multi-agent joint evolution, which will offer novel insights to the community. Specifically, we first conduct node-wise Monte Carlo sampling, which samples the behaviors of each agent across the entire reasoning chain during joint evolution to stabilize reward estimation and reduce computational overhead. Then, the value discounting is applied to the current agent action by integrating decision efficiency to perform end-to-end scoring. However, updating all agents simultaneously is also challenging and remains unsolved. Here, we follow the marginal benefit principle [31] to update agents associated with the top-$K$ sampling groups exhibiting the largest variance of intra-group rewards, maximizing joint utility at minimal computational cost. To facilitate the multi-agent training, we design a memory evolution mechanism harnessed from GRPO rewards as a "free lunch" since our key observation is that these rewards are inherently coupled to memory. Through direct utilization of GRPO rewards to update reasoning-associated memory, the decision-making and memory modules achieve synchronous optimization, effectively alleviating the difficulty of joint evolution.

To sum up, the main contributions of this work can be listed as follows:

- We investigate an effective methodology for integrating GRPO into heterogeneous multi-agent systems and introduce JoyAgents-R1, a joint evolution dynamics that enables each agent to improve its decision-making and memory capabilities, achieving an optimal balance with superior overall performance.

- We propose a novel strategy that combines node-wise Monte Carlo sampling with the marginal benefit principle to address the combinatorial explosion of trajectories during multi-agent GRPO updates, thereby enhancing sampling efficiency and training stability.

- We design a simple yet effective memory mechanism derived from GRPO rewards to facilitate the multi-agent evolution. During LLM weight updates, agents' memories evolve based on action rewards, significantly accelerating convergence and boosting reasoning performance.

- We conduct extensive experiments on both general and vertical domain benchmarks to validate the effectiveness of our method. JoyAgents-R1 could generate more reasonable and accurate results based on smaller open-source models.

## 2 Related Work

### 2.1 LLM-based multi-agent planning

Recent breakthroughs in LLMs have transformed the landscape of agent planning [32]. Autonomous agents can implement iterative self-reflection mechanisms, dynamically integrate external information via structured prompts [33–37], and perceive environments to plan tasks through sophisticated reasoning and decision-making processes [38]. Compared with single-agent approaches that struggle with inefficiency and environmental adaptability, multi-agent systems achieve robust performance through decentralized decision-making and collaborative mechanisms, enabling the coordination of agents with distinct capabilities and objectives to pursue shared goals in fields like robotics [39], tool calling [40], and AI assistants [11]. However, few multi-agent systems can achieve multi-domain tasks. Moreover, open-source LLMs [3, 5] lag significantly behind state-of-the-art models, which are either closed-source with opaque mechanisms [1, 2] or overly complex for multi-agent deployment [23]. This work introduces a hierarchical multi-agent architecture to interpret user

queries and perform dynamic planning. Based on the smaller open-source LLMs, our framework implements diverse capabilities, including question answering, mathematical computation, and tool calling, revealing the mechanisms that drive effective heterogeneous multi-agent collaboration in resource-constrained environments.

## 2.2 Multi-agent reinforcement learning

MARL has witnessed substantial advancements, rendering it an ideal approach for tackling complex and challenging tasks [41]. This work focuses on cooperative MARL tasks where various agents share a common goal, which has been successfully applied in many fields such as game playing [20, 42], task allocation [15], skill discovery [43], and circuit design [44]. Typical MARL methods employ an actor-critic framework, where actors generate actions based on observations, and critics evaluate their long-term efficacy [45]. There are policy-based methods like MADDPG [46] and MAPPO [20] and value-based ones like VDN [47] and QMIX [48]. Although recent studies have explored LLM-based MARL frameworks for problem-solving [49, 50] and embodied intelligence [39, 51, 52], these approaches primarily focus on enhancing inter-agent communication and cooperative decision-making, with limited emphasis on the joint evolution of multi-agent systems. In addition, many methods adopt parameter sharing across agents, which restricts their applicability to homogeneous scenarios [25, 53, 54] and fails to address heterogeneous systems [22]. Recently, GRPO [24] has gained popularity through DeepSeek-R1 [23]. The algorithm eliminates the value function and relies on the observed rewards, making it well-suited for joint training of complicated heterogeneous multi-agents [55]. However, due to the dynamic and changeable reasoning paths in multi-agents, direct sampling based on GRPO will lead to an exponential explosion. Therefore, we propose a novel strategy that leverages node-wise Monte Carlo with marginal benefit optimization to guide reasoning path sampling and agent updating, improving sampling efficiency and stabilizing training dynamics.

## 2.3 LLM-based agent memory

Agent memory can be divided into RAG-based and embodied categories [32]. The former is typically stored in additional storage, while the latter embeds memories into model parameters by fine-tuning LLMs. In this work, we focus on RAG-based long-term memory mechanisms. Recent works have explored diverse strategies [56–59]. Methods such as MoT [60], TiM [61], and RAP [62] aim to improve LLM reasoning and planning by leveraging memories after selection or thinking. MemoryBank [63] draws inspiration from the Ebbinghaus forgetting curve [64] to design a selective information retention mechanism. HELPER [65], ExpeL [66], RET-LLM [67], Synapse [68], and A-mem [69] adopt different approaches for knowledge aggregation, storage, and retrieval, enhancing LLMs' adaptability to novel tasks. Moreover, there are several memory mechanisms tailored for multi-agent systems [70], exploring the memory synchronization [71], communication [51], and the information asymmetry [72] among agents. Nevertheless, existing memory modules often struggle to synchronize with LLM updates, limiting system efficacy. In contrast, we we present the joint evolution dynamics where agent memory and decision-making modules evolve synergistically with LLMs optimization. This mechanism leverages GRPO rewards as cost-free supervisory signals, eliminating the need for dedicated model training while enhancing convergence efficiency.

## 3 Method

We introduce JoyAgents-R1, a novel joint evolution dynamics for multi-agent reinforcement learning. First, a hierarchical architecture is designed to integrate heterogeneous multi-agents for complex collaborative tasks (Section 3.1). Then, a variance-reduction GRPO including node-wise Monte Carlo sampling and marginal benefit-driven updating is constructed for joint training (Section 3.2). Finally, an adaptive memory evolution mechanism leveraging GRPO rewards as cost-free supervisory signals is proposed to enable synchronized optimization of agent decision modules (Section 3.3).

## 3.1 The architecture of JoyAgents-R1

As illustrated in Fig. 1, the proposed JoyAgents-R1 adopts a hierarchical architecture, consisting of a master agent and multiple sub-agents as follows:
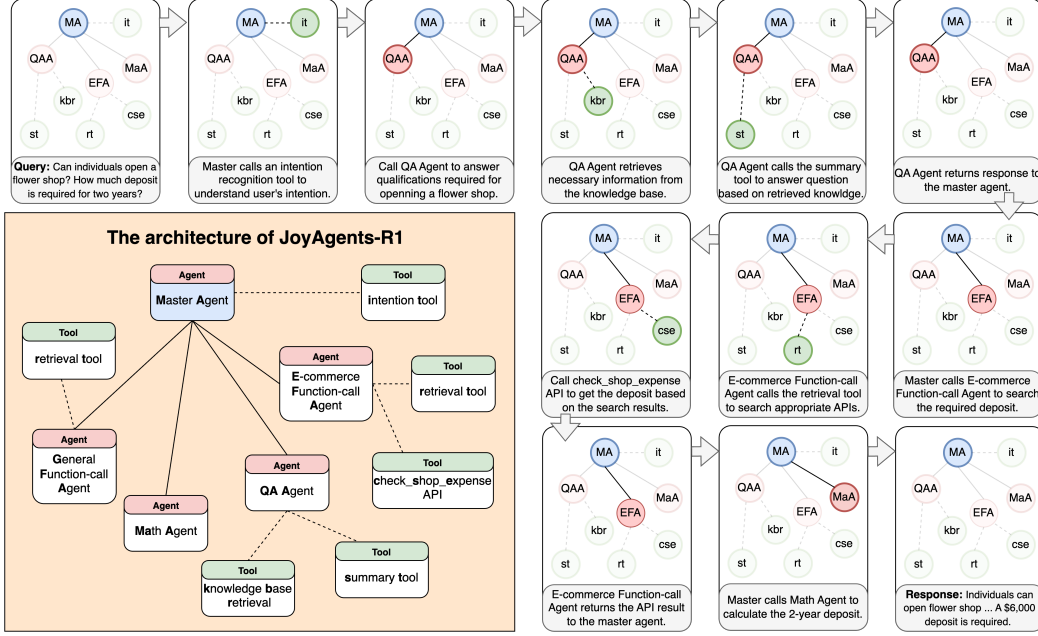
Figure 1: The multi-agent architecture and a reasoning example of JoyAgents-R1.

- **Master agent** first analyzes the query, then orchestrates sub-agents (*e.g.*, question answering) or tools (*e.g.*, intention recognition) in each step, and determines the final response to the user.
- **Question-answering agent** could perform general and domain-specific (*e.g.*, e-commerce) question answering by retrieving and summarizing the recalled information from external knowledge bases.
- **Function-call agents** include general-purpose and domain-specific (*e.g.*, e-commerce) agents. These agents either execute the function call directly via memory-driven prompts or invoke a tool retriever and make further selections from the recalled APIs.
- **Math agent** is specialized for solving mathematical problems.

Fig. 1 shows an example of a reasoning chain during inference. Upon receiving a query, the master agent analyzes the user intent and assigns tasks to specialized sub-agents, which execute iterative operations using contextual inputs transmitted by the master until completion. Then, the results are relayed back to the master for subsequent planning cycles until termination. Each agent executes in a ReAct manner [33] and dynamically retrieves validated strategies from memory to minimize redundant reasoning and enhance decision efficiency. On the other hand, since the inherent complexity of actor-critic frameworks with critic models, JoyAgents-R1 employs Group Relative Policy Optimization (GRPO) for policy optimization during training, which foregoes the value function and computes advantages in a group-relative manner [24]. However, its direct application in multi-agent joint evolution faces the following challenges:

- **Low sampling efficiency.** Since multiple actions are sampled for each agent throughout the reasoning chain, the number of trajectories explodes exponentially.
- **Inefficient parameter updates.** Sequential policy updates for all agents in the chain will lead to inefficient parameter optimization. Furthermore, decoupled policy adjustments struggle to coordinate inter-agent dependencies for overall performance enhancement.
- **Slow training convergence.** Since diverse architectures and output impedes policy synchronization, heterogeneous multi-agent systems exhibit pronounced convergence difficulties during training.

### 3.2 Variance reduction GRPO tailored for multi-agent system

To address the above challenges, the joint evolution dynamics, which introduces a variance-reducing GRPO tailored for multi-agent is proposed (Fig. 2). The method consists of the following parts:
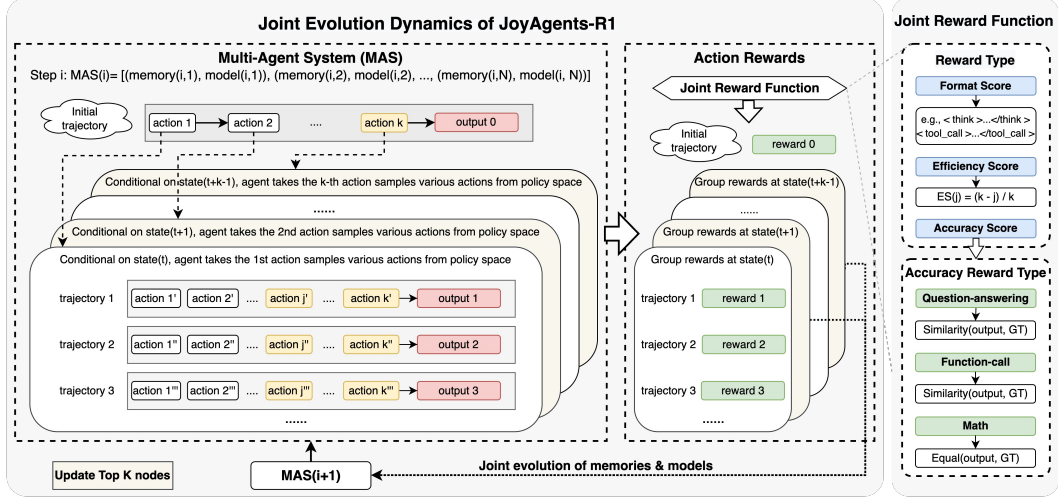
Figure 2: Joint evolution dynamics of JoyAgents-R1.

**Node-wise Monte Carlo sampling.** For a multi-agent system (MAS) consisting of $N$ agents $\{m_i\}_{i=1}^{N}$, each agent $m_i$ has an action space including $G_i$ possible actions. Therefore, when GRPO is directly applied to MAS in a naive way, a reasoning trajectory of length $k$ will generate $G_{Naive} = G_1 \times G_2 \cdots \times G_k$ sampling paths in an exponential explosion. Different from that, JoyAgents-R1 performs node-wise Monte Carlo sampling as shown in Fig. 2. Specifically, given a query $q$ obtained from the dataset, JoyAgents-R1 first produces an initial trajectory of length $k$, and then sequentially samples $(G_i - 1)$ actions for each node $m_i$ in the trajectory. To facilitate comparative calculations and solve the exponential explosion problem, the reasoning path from the query to the sampling node (*e.g.*, $m_j$) still follows the original trajectory. After $m_j$, no more sampling is conducted and reasoning continues until the end. In this way, a total of $G_{NMC} = G_1 + G_2 \cdots + G_k$ trajectories are generated. This addition operation is much smaller than the multiplication one, which effectively improves the sampling efficiency, especially for MAS with long trajectories.

**Model update based on marginal benefit principle.** For the $k$ groups of sampled trajectories, the model parameters of corresponding agents $\{m_i\}_{i=1}^{k}$ are $\Theta = \{\theta_i\}_{i=1}^{k}$. Therefore, the objective for updating all models in a straightforward manner could be calculated as follows:

$$\mathcal{J}_{Multi-GRPO}(\Theta) = \{\mathcal{J}_{GRPO}(\theta_i) \mid i = 1, \cdots, k\} \tag{1}$$

where $\mathcal{J}_{GRPO}(\theta_i)$ is objective for the policy model $\pi_{\theta_i}$ from DeepSeek-R1 [23]. In addition, the original KL penalty between current and reference policies is eliminated like [73], allowing greater distributional divergence from the initial model while eliminating the need for multiple reference models. This is more suitable for our multi-agent coordination task, where strategy heterogeneity and dynamic adaptation are crucial.

Furthermore, to accelerate the policy update and enable the perception of global states, this method implements a variance reduction GRPO based on the marginal benefit principle [31] as follows:

$$\mathcal{J}_{VR-GRPO}(\Theta) = \{\mathcal{J}_{GRPO}(\theta_i) \mid i \in \text{argtopK}\left(\text{Var}(R_i), K\right)\} \tag{2}$$

where $R_i = \{r_j\}_{j=1}^{G_i}$ is the reward set obtained after sampling $G_i$ actions from $m_i$ to the end of trajectory. argtopK $(\cdot)$ returns the index set of the first $K$ nodes with the largest reward variance. To this end, JoyAgents-R1 prioritizes updating model parameters for the top-$K$ agents exhibiting the largest performance fluctuations among all reasoning trajectory participants. Compared to updating all agents sequentially as shown in Equation 1, this variance-aware selection strategy minimizes computational overhead while maximizing the joint benefit, efficiently steering multi-agent parameter updates through GRPO rewards in a dynamic paradigm.

**Action rewards.** As shown in Fig. 2, given the final answer $a$ corresponding to $q$ as the ground truth (GT), the reward $\mathcal{R}$ (*i.e.*, $r_i$ in the above text) of each agent action consists of three terms as follows:

$$\mathcal{R} = \mathcal{R}_{\mathcal{A}} + \mathcal{R}_{\mathcal{F}} - \mathcal{R}_{\mathcal{E}} \tag{3}$$

5

**Algorithm 1** Dynamic Memory Updating from GRPO Rewards

---

**Input:** Planning chain $\mathcal{P}$, Trajectory output $\mathcal{O}$, Trajectory reward $\mathcal{R}_\mathcal{M}$, Upper bound $U$, Lower bound $L$,
       Set of $n$ recalled memories $Recall_q$ given query $q$, Hyperparameters $\alpha$ and $\beta$, Timestamp $t$

1: **if** $\mathcal{R}_\mathcal{M} > U$ **then**
2:     Insert new memory: $\mathcal{M}_{new} = (q, t, \mathcal{P}, \mathcal{O}, \mathcal{R}_\mathcal{M})$
3: **end if**
4: **for** each memory $\mathcal{M}_i \in Recall_q$ **do**
5:     Compute similarity: $sim_i \leftarrow Sim(\mathcal{O}, \mathcal{O}_i)$ (Direct answer mode) **or** $Sim(\mathcal{P}, \mathcal{P}_i)$ (Tool call mode)
6: **end for**
7: Normalize similarity scores: $s_i \leftarrow \text{Softmax}(\{sim_j\}_{j=1}^n)$
8: **for** each memory $\mathcal{M}_i \in Recall_q$ **do**
9:     **if** $\mathcal{R}_\mathcal{M} > U$ **then**
10:        Update timestamp: $t_i \leftarrow t$;
11:        Compute time and reward differences: $\Delta t \leftarrow -|t - t_i|$ **and** $\Delta s \leftarrow s_i \cdot |\mathcal{R}_\mathcal{M} - U|$
12:     **else if** $\mathcal{R}_\mathcal{M} < L$ **then**
13:        Compute time and reward differences: $\Delta t \leftarrow -|t - t_i|$ **and** $\Delta s \leftarrow -s_i \cdot |\mathcal{R}_\mathcal{M} - L|$
14:     **end if**
15:     Update recalled memory scores: $\mathcal{R}_{\mathcal{M}i} \leftarrow \mathcal{R}_{\mathcal{M}i} + \alpha\Delta t + \beta\Delta s$
16: **end for**

---

- *Accuracy reward ($\mathcal{R}_\mathcal{A}$).* The accuracy reward is calculated end-to-end by using the trajectory final answer. The evaluation methods are tailored to tasks' distinct settings and output formats. For instance, semantic similarity metrics assess alignment with ground truth for tasks like question answering and function calling, while mathematical operations require exact-match validation against predefined solutions.

- *Format reward ($\mathcal{R}_\mathcal{F}$).* Since each agent infers in the RaAct style, the model output is formatted using HTML tags for thinking (*i.e.*, $< think > \cdots < /think >$ ) and tool calling (*i.e.*, $< tool\_call > \cdots < /tool\_call >$ ). Format rewards could guide the model to generate structured results, improving clarity and enhancing the reasoning ability of LLMs.

- *Efficiency reward ($\mathcal{R}_\mathcal{E}$).* For the $j$-th node in a trajectory of length $k$, its efficiency score is computed as $\mathcal{R}_\mathcal{E} = \frac{k-j}{k}$, which imposes a penalty proportional to the distance of the node from the trajectory endpoint, where the efficiency score is quantified by number of subsequent decision steps required, showing how the current node's plan impacts downstream computational cost.

## 3.3   Free lunch in GRPO rewards for memory evolution

For the heterogeneous MAS, the memory modules are introduced to accelerate the model training and reduce unnecessary repeated reasoning. To improve the efficiency of joint evolution, the agent memory is designed to undergo dynamic adaptation alongside LLMs' update during training (Fig. 2). Different from previous methods that require training dedicated models or utilizing LLMs to evolve memory, JoyAgents-R1 creates a simple yet effective memory updating mechanism that leverages GRPO rewards as a cost-free supervisory signal and mainly consists of the following three steps:

**Adaptive reward thresholding.** The memory of each agent along a trajectory is updated using a unified reward without efficiency score (*i.e.*, $\mathcal{R}_\mathcal{M} = \mathcal{R}_\mathcal{A} + \mathcal{R}_\mathcal{F}$), which ensures consistent update criteria and enables independent memory modules to perceive the overall performance. For $G_{NMC}$ trajectories sampled from query $q$, the mean $\mu$ and standard deviation $\sigma$ of corresponding rewards are first computed, then the 2.5% and 97.5% percentiles of the approximate normal distribution are selected as the lower ($L = \mu - 1.96\sigma$) and upper ($U = \mu + 1.96\sigma$) bounds for memory updates.

**Dynamic memory updating.** As shown in Algorithm 1, the algorithm determines whether to add a new memory based on the reward threshold and dynamically updates the memory recalled by the query. In addition, other memories are updated only according to time decay: $\mathcal{R}_{\mathcal{M}i} \leftarrow \mathcal{R}_{\mathcal{M}i} + \alpha\Delta t$.

**Memory overflow handling.** To ensure memory quality and save storage space, the memory will be deleted either when its final reward ($\mathcal{R}_{\mathcal{M}i}$) falls below a predefined threshold $D$ or when the memory capacity exceeds upper bounds and the memory's reward rank is relatively low.

To this end, the memory module is synergistically updated with model parameters through trajectory rewards, accelerating training convergence and boosting reasoning performance.

# 4 Experiments

## 4.1 Implementation details

In this work, we opt for Qwen2.5-3B [5] as the backbone for each agent to ensure technical reproducibility. The experiment consists of two main stages. In the first stage, the base models are fine-tuned with a learning rate of 5e-6 for 5 epochs. In the second stage, the multi-agent system is trained via reinforcement learning at a learning rate of 1e-6 for 5 epochs. Specifically, each agent from an initial trajectory is sampled $G_i = 5$ actions with a temperature of 1.2. Subsequently, the top-5 nodes are selected for model updates. Similar to DeepSeek-R1 [23], the iterative RL with GRPO is executed for 2 iterations. Regarding memory evolution, the deletion threshold is set to $D = 0$, with hyperparameters $\alpha = \beta = 1$. The models are trained on 8 NVIDIA H200 GPUs, and the best results are reported. More details are provided in the supplementary material (Section A).

## 4.2 Datasets and setup

To verify the effectiveness of the proposed method, we construct a heterogeneous multi-agent dataset for AI assistant scenarios, including general and e-commerce fields as follows:

**Supervised fine-tuning datasets.** A dedicated dataset is constructed for each agent. The input integrates diverse elements, such as user queries, retrieved memory, optional tools, historical dialogues, and tool-generated responses. The target comprises the reasoning processes (*i.e.*, $< think > \cdots < /think >$), tool calling (*i.e.*, $< tool\_call > \cdots < /tool\_call >$), or final answers. Specifically, the master agent, tasked with dynamic reasoning and orchestrating four sub-agents, uses a training set of 13,000 samples: 10,000 individual sub-agent calls (2,500 per sub-agent) and 3,000 collaborative agent invocations. The QA agent is trained on 1,000 instances, comprising 700 real-world e-commerce data entries and 300 open-domain cases from COIG [74]. Regarding function-calling, the e-commerce agent encompasses 12 prevalent APIs of e-commerce platforms with 2,000 training cases, while the general-purpose agent leverages 1,000 diverse API calls from ToolBench [9], yielding 3,500 training instances. For the math agent, GSM8K [75] is used to construct the dataset.

**Reinforcement learning datasets.** The RL dataset contains query-response pairs with ground-truth annotations. The training set includes 100 samples per sub-agent task (*i.e.*, math, e-commerce/general function calls, QA) and 200 collaborative instances. The test set has 500 samples, including 100 instances for each independent task and 100 cases for the collaborative task.

## 4.3 Evaluation metrics

The experimental evaluation metrics are primarily categorized into accuracy and efficiency. For accuracy, the calculation methods vary across task types. In question-answering tasks, accuracy is measured by the semantic similarity between the predicted and ground truth answers, with a threshold of 0.6. For mathematical problems, accuracy is binary (0 or 1) based on exact numerical matching. For function-call tasks, a response with correct function names is scored 1, otherwise 0. Regarding efficiency, it is quantified by the average number of steps required to complete a task, reflecting the reasoning efficiency of the method.

## 4.4 Comparative experiments

**Comparison rules.** To verify the effectiveness of the proposed method, a comprehensive comparison is conducted with agents based on closed-source and open-source models. All baselines run under a consistent set of prompts, and each model is assessed as a single agent through react-based multi-turn interactions. Note that the Qwen2.5-7B is trained by using SFT, while other baselines use native large models. For the proposed multi-agent method, the model parameters activated for each query typically range from 6B (master plus a single sub-agent) to a maximum of 15B (master plus 4 sub-agents).

**Comparison with closed-source models.** The upper section of Table 1 presents a comparative analysis with closed-source models, which generally exhibit a superior average accuracy metric compared to our approach. Given that these models represent the current state-of-the-art in general-purpose language understanding, surpassing DeepSeek-V3 in average accuracy metric constitutes a significant achievement for our specialized framework. A particularly compelling observation emerges from the E-commerce Function-Call subtask evaluation: JoyAgents-R1 achieves 24% higher

Table 1: Accuracies (%) of multi-tasks with agents based on larger SOTA closed-source or open-source models. 'FC' is the function call. '*' indicates using the native model. '†' is the best score.

| Model | Math | QA | E-commerce FC | General FC | Cooperation | Average |
|---|---|---|---|---|---|---|
| DeepSeek-R1 * | 98.0† | 37.0 | 24.0 | 72.0 | 7.0† | 47.6 |
| DeepSeek-V3 * | 96.0 | 32.0 | 10.0 | 68.0 | 2.0 | 41.6 |
| GPT-4o * | 85.0 | 35.0 | 56.0† | 83.0† | 6.0 | 53.0† |
| Qwen2.5-32B * | 72.0 | **38.0†** | 3.0 | 72.0 | 2.0 | 37.4 |
| Qwen2.5-14B * | **80.0** | 29.0 | 0.0 | 42.0 | 1.0 | 30.4 |
| Qwen2.5-7B (SFT) | 65.0 | 13.0 | 42.0 | 73.0 | 3.0 | 39.2 |
| **JoyAgents-R1** (3B) | 68.0 | 22.0 | **48.0** | **76.0** | **6.0** | **44.0** |

Table 2: Ablation study for multi-agent reasoning on multiple tasks. From left to right, whether to train with RL (*vs.* SFT), whether to generate the think process, whether to update top-$K$ models (*vs.* update all models), whether to use efficiency rewards, and whether to integrate memory modules.

| Method | RL | Think | Top-$K$ | Efficiency | Memory | Accuracy↑ | Reasoning Rounds↓ |
|---|---|---|---|---|---|---|---|
| $M_1$ | ✗ | ✗ | - | - | - | 17.2 | **5.4** |
| $M_2$ | ✗ | ✔ | - | - | - | 35.0 | 7.1 |
| $M_3$ | ✔ | ✔ | ✗ | ✔ | ✔ | 40.0 | 7.4 |
| $M_4$ | ✔ | ✔ | ✔ | ✗ | ✔ | 42.4 | 8.1 |
| $M_5$ | ✔ | ✔ | ✔ | ✔ | ✗ | 40.0 | 7.7 |
| $M_6$ | ✔ | ✔ | ✔ | ✔ | ✔ | **44.0** | 7.8 |

accuracy than DeepSeek-R1, and outperforms DeepSeek-V3 by 38%, while trailing GPT-4o by only 8%. We also achieved performance comparable to state-of-the-art methods in the General FC subtask. This significant difference highlights two key insights: (1) Conventional SOTA models exhibit limited generalizability when directly applied to specialized vertical domains without domain adaptation. (2) Compact architectures like ours can achieve competitive performance through meticulous domain-specific optimization. These findings reveal the viability of targeted architectures and training strategies for domain-specific applications.

**Comparison with open-source models.** The lower section of Table 1 presents a comparative analysis with open-source models, where we selected the Qwen-2.5 series models (7B, 14B, and 32B) to maintain consistency with the base model architecture. Notably, the 7B variant is a fine-tuned version. Empirical results demonstrate that our method achieves best performance, surpassing all open-source counterparts by margins of 4.8%, 13.6%, and 6.6% respectively. This superiority is particularly pronounced in the E-commerce FC subtask, where our approach outperforms even the fine-tuned 7B model, underscoring the efficacy of our reinforcement learning framework. The relatively lower performance on the math subtask compared to other models can be attributed to the disparity in model scale. The consistent outperformance of our method over other larger-scale models not only validates the effectiveness of our methodology but also highlights its parameter efficiency in resource-constrained scenarios.

### 4.5 Ablation study

Table 2 illustrates the quantitative ablation results for multi-agent reasoning as follows:

**Effectiveness of generating the think process.** $M_1$ separately trains each agent's base model via SFT, yielding the lowest accuracy (17.2%). Compared to $M_1$, $M_2$ further incorporates think process generation and significantly improves accuracy to 35.0%, demonstrating the indispensable role of explicit reasoning in agent decision-making.

**Effectiveness of training with reinforcement learning.** Different from $M_2$, $M_6$ integrates GRPO for the joint evolution of multi-agents and enhances the accuracy by 25.7% relatively. Furthermore, the results from $M_3$ to $M_6$ using RL are much better than $M_2$, indicating the effectiveness of global perception and preference alignment in joint training with reinforcement learning.
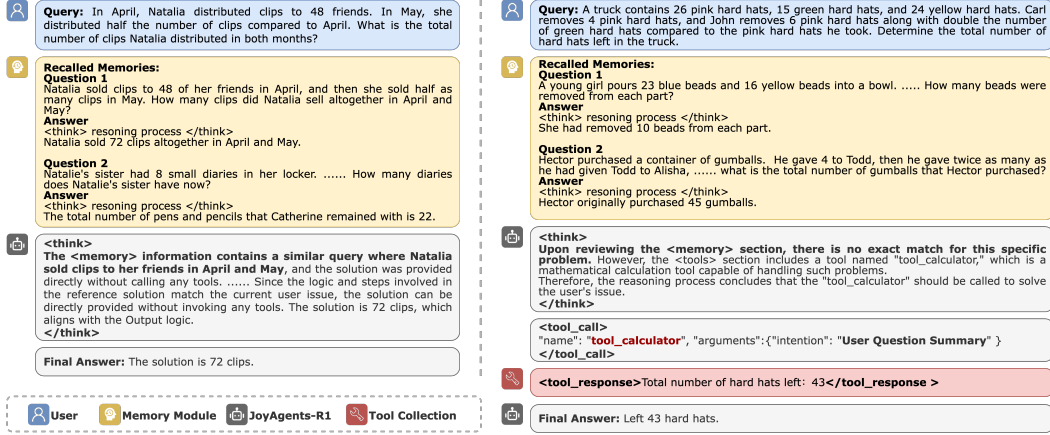
Figure 3: Examples of different agents making decisions based on dynamically recalled memories.

**Effectiveness of updating top-$K$ models.** $M_3$ sequentially updates all models (Eq. 1), while $M_6$ improves its accuracy by 10% relatively through targeted updates of top-k nodes with maximal reward fluctuations. Inspired by marginal benefit, $M_6$ prioritizes models that need the most optimization in the global state to achieve the maximum benefit with the minimum cost and improve performance.

**Effectiveness of utilizing efficiency rewards.** Since $M_4$ excludes $\mathcal{R}_\mathcal{E}$ in Eq. 3 during training, it exhibits suboptimal accuracy (42.4%) with a maximum reasoning round of 8.1, demonstrating that efficiency constraints are critical for balancing performance and computational cost. In contrast, $M_1$ consumes the fewest inference rounds due to insufficient reasoning capacity, which forces premature termination from ineffective decision-making.

**Effectiveness of integrating memory modules.** $M_5$ disables memory modules across all agents, forcing reliance on decision modules for planning and reasoning. This leads to a relative decrease of 10.0% in accuracy compared to $M_6$, indicating the important role of memory modules in enhancing model performance through joint evolution and reasoning enhancement.

### 4.6 Case analysis

Fig. 3 shows that JoyAgents-R1 dynamically retrieves memories through query similarity to guide decision-making. On the left, JoyAgents-R1 identifies recalled memories containing a question that is semantically similar and numerically identical to the user's query. So it directly reuses stored answers to avoid duplicate reasoning and enhance response efficiency. For the right case, when there is no useful instance from recalled memories, our method deliberates and opts to invoke a specialized math tool, enabling problem-solving and accurate resolution. These cases validate JoyAgents-R1's flexible integration and application of decision-making and memory modules to solve complicated tasks.

## 5 Conclusion and Discussion

**Conclusion.** This paper introduces JoyAgents-R1, a joint evolution dynamics for heterogeneous multi-agents integrating GRPO. First, the node-wise Monte Carlo sampling and marginal benefit-driven evolution strategies are implemented to enhance sampling efficiency and training stability. Then, an adaptive mechanism using GRPO rewards is designed to guide memory updating, alleviating redundant reasoning and promoting convergence. Experimental results demonstrate that our method achieves comparable performance to larger LLMs, even based on smaller open-source models.

**Limitations.** Due to computational constraints, the proposed method currently focuses on multi-agent systems based on small open-source models. In addition, existing training frameworks lack compatibility with heterogeneous multi-agent joint training. Future work will focus on scaling up LLMs to achieve performance gains and cross-domain robustness, while engineering a dedicated framework optimized for heterogeneous multi-agent co-evolution with computational efficiency.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[2] Anthropic. The claude 3 model family: Opus, sonnet, haiku. 2024.

[3] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

[4] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

[5] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

[6] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. In *IJCAI*, 2024.

[7] Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. Api-bank: A comprehensive benchmark for tool-augmented llms. *arXiv preprint arXiv:2304.08244*, 2023.

[8] Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu, Tianpeng Bao, Guoqing Du, Shiwei Shi, Hangyu Mao, Ziyue Li, Xingyu Zeng, et al. Tptu: Large language model-based ai agents for task planning and tool usage. *arXiv preprint arXiv:2308.03427*, 2023.

[9] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. In *ICLR*, 2024.

[10] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, et al. A survey on in-context learning. In *EMNLP*, 2024.

[11] Yicheng Fu, Raviteja Anantha, and Jianpeng Cheng. Camphor: Collaborative agents for multi-input planning and high-order reasoning on device. *arXiv preprint arXiv:2410.09407*, 2024.

[12] Amutheezan Sivagnanam, Ava Pettet, Hunter Lee, Ayan Mukhopadhyay, Abhishek Dubey, and Aron Laszka. Multi-agent reinforcement learning with hierarchical coordination for emergency responder stationing. *arXiv preprint arXiv:2405.13205*, 2024.

[13] Zepeng Ning and Lihua Xie. A survey on multi-agent reinforcement learning and its application. *JAI*, 2024.

[14] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022.

[15] Shariq Iqbal, Robby Costales, and Fei Sha. Alma: Hierarchical learning for composite multi-agent tasks. In *NeurIPS*, 2022.

[16] Zikang Tian, Ruizhi Chen, Xing Hu, Ling Li, Rui Zhang, Fan Wu, Shaohui Peng, Jiaming Guo, Zidong Du, Qi Guo, et al. Decompose a task into generalizable subtasks in multi-agent reinforcement learning. In *NeurIPS*, 2023.

[17] Wei Fu, Chao Yu, Zelai Xu, Jiaqi Yang, and Yi Wu. Revisiting some common practices in cooperative multi-agent reinforcement learning. In *ICML*, 2022.

[18] Huao Li, Hossein Nourkhiz Mahjoub, Behdad Chalaki, Vaishnav Tadiparthi, Kwonjoon Lee, Ehsan Moradi Pari, Charles Lewis, and Katia Sycara. Language grounded multi-agent reinforcement learning with human-interpretable communication. In *NeurIPS*, 2024.

[19] Pablo Hernandez-Leal, Michael Kaisers, Tim Baarslag, and Enrique Munoz De Cote. A survey of learning in multiagent environments: Dealing with non-stationarity. *arXiv preprint arXiv:1707.09183*, 2017.

[20] Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative multi-agent games. In *NeurIPS*, 2022.

[21] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[22] Yifan Zhong, Jakub Grudzien Kuba, Xidong Feng, Siyi Hu, Jiaming Ji, and Yaodong Yang. Heterogeneous-agent reinforcement learning. *JMLR*, 2024.

[23] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

[24] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

[25] Yue Deng, Weiyu Ma, Yuxin Fan, Ruyi Song, Yin Zhang, Haifeng Zhang, and Jian Zhao. Smac-r1: The emergence of intelligence in decision-making tasks. *arXiv preprint arXiv:2410.16024*, 2024.

[26] Xiaobo Xia and Run Luo. Gui-r1 : A generalist r1-style vision-language action model for gui agents. *arXiv preprint arXiv:2504.10458*, 2025.

[27] Zhengxi Lu, Yuxiang Chai, Yaxuan Guo, Xi Yin, Liang Liu, Hao Wang, Guanjing Xiong, and Hongsheng Li. Ui-r1: Enhancing action prediction of gui agents by reinforcement learning. *arXiv preprint arXiv:2503.21620*, 2025.

[28] Wenxuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Yao Hu, and Shaohui Lin. Vision-r1: Incentivizing reasoning capability in multimodal large language models. *arXiv preprint arXiv:2503.06749*, 2025.

[29] Pablo Hernandez-Leal, Bilal Kartal, and Matthew E Taylor. A very condensed survey and critique of multiagent deep reinforcement learning. In *AAMAS*, 2020.

[30] Qihan Liu, Jianing Ye, Xiaoteng Ma, Jun Yang, Bin Liang, and Chongjie Zhang. Efficient multi-agent reinforcement learning by planning. In *ICLR*, 2024.

[31] Emil Kauder. *History of marginal utility theory*. Princeton University Press, 2015.

[32] Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. Understanding the planning of llm agents: A survey. *arXiv preprint arXiv:2402.02716*, 2024.

[33] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *ICLR*, 2023.

[34] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. In *NeurIPS*, 2023.

[35] Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language agent tree search unifies reasoning, acting, and planning in language models. In *ICML*, 2024.

[36] Bill Yuchen Lin, Yicheng Fu, Karina Yang, Faeze Brahman, Shiyu Huang, Chandra Bhagavatula, Prithviraj Ammanabrolu, Yejin Choi, and Xiang Ren. Swiftsage: A generative agent with fast and slow thinking for complex interactive tasks. In *NeurIPS*, 2023.

[37] Shuofei Qiao, Ningyu Zhang, Runnan Fang, Yujie Luo, Wangchunshu Zhou, Yuchen Eleanor Jiang, Huajun Chen, et al. Autoact: Automatic agent learning from scratch for qa via self-planning. In *ICLR Workshop*, 2024.

[38] Malik Ghallab, Dana Nau, and Paolo Traverso. *Automated Planning: theory and practice*. Elsevier, 2004.

[39] Shyam Sundar Kannan, Vishnunandan LN Venkatesh, and Byung-Cheol Min. Smart-llm: Smart multi-agent robot task planning using large language models. In *IROS*, 2024.

[40] Weizhou Shen, Chenliang Li, Hongzhan Chen, Ming Yan, Xiaojun Quan, Hehong Chen, Ji Zhang, and Fei Huang. Small llms are weak tool learners: A multi-llm agent. In *EMNLP*, 2024.

[41] Lei Yuan, Yuqi Bian, Lihe Li, Ziqian Zhang, Cong Guan, and Yang Yu. Efficient multi-agent offline coordination via diffusion-based trajectory stitching. In *ICLR*, 2025.

[42] T Wang, T Gupta, B Peng, A Mahajan, S Whiteson, and C Zhang. Rode: learning roles to decompose multi- agent tasks. In *ICLR*, 2021.

[43] Mingyu Yang, Yaodong Yang, Zhenbo Lu, Wengang Zhou, and Houqiang Li. Hierarchical multi-agent skill discovery. In *NeurIPS*, 2023.

[44] Zhihai Wang, Jie Wang, Dongsheng Zuo, Ji Yunjie, Xilin Xia, Yuzhe Ma, Jianye Hao, Mingxuan Yuan, Yongdong Zhang, and Feng Wu. A hierarchical adaptive multi-task reinforcement learning framework for multiplier circuit design. In *ICML*, 2024.

[45] Chuanneng Sun, Songjun Huang, and Dario Pompili. Llm-based multi-agent reinforcement learning: Current and future directions. *arXiv preprint arXiv:2405.11106*, 2024.

[46] Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30, 2017.

[47] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *AAMAS*, 2018.

[48] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 2020.

[49] Huaben Chen, Wenkang Ji, Lufeng Xu, and Shiyu Zhao. Multi-agent consensus seeking via large language models. *arXiv preprint arXiv:2310.20151*, 2023.

[50] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. Metagpt: Meta programming for a multi-agent collaborative framework. In *ICLR*, 2024.

[51] Zhao Mandi, Shreeya Jain, and Shuran Song. Roco: Dialectic multi-robot collaboration with large language models. In *ICRA*, 2024.

[52] Hongxin Zhang, Weihua Du, Jiaming Shan, Qinhong Zhou, Yilun Du, Joshua B Tenenbaum, Tianmin Shu, and Chuang Gan. Building cooperative embodied agents modularly with large language models. In *ICLR*, 2024.

[53] Nathaniel Grammel, Sanghyun Son, Benjamin Black, and Aakriti Agrawal. Revisiting parameter sharing in multi-agent deep reinforcement learning. *arXiv preprint arXiv:2005.13625*, 2020.

[54] Filippos Christianos, Georgios Papoudakis, Muhammad A Rahman, and Stefano V Albrecht. Scaling multi-agent reinforcement learning with selective parameter sharing. In *ICML*, 2021.

[55] Zixuan Ke, Fangkai Jiao, Yifei Ming, Xuan-Phi Nguyen, Austin Xu, Do Xuan Long, Minzhi Li, Chengwei Qin, Peifeng Wang, Silvio Savarese, et al. A survey of frontiers in llm reasoning: Inference scaling, learning to reason, and agentic systems. *arXiv preprint arXiv:2504.09037*, 2025.

[56] Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G Patil, Ion Stoica, and Joseph E Gonzalez. Memgpt: Towards llms as operating systems. *arXiv preprint arXiv:2310.08560*, 2023.

[57] Junru Lu, Siyu An, Mingbao Lin, Gabriele Pergola, Yulan He, Di Yin, Xing Sun, and Yunsheng Wu. Memochat: Tuning llms to use memos for consistent long-range open-domain conversation. *arXiv preprint arXiv:2308.08239*, 2023.

[58] Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu, Xifeng Yan, Jianfeng Gao, and Furu Wei. Augmenting language models with long-term memory. 2023.

[59] Ziheng Huang, Sebastian Gutierrez, Hemanth Kamana, and Stephen MacNeil. Memory sandbox: Transparent and interactive memory management for conversational agents. In *UIST*, 2023.

[60] Xiaonan Li and Xipeng Qiu. Mot: Memory-of-thought enables chatgpt to self-improve. In *EMNLP*, 2023.

[61] Lei Liu, Xiaoyan Yang, Yue Shen, Binbin Hu, Zhiqiang Zhang, Jinjie Gu, and Guannan Zhang. Think-in-memory: Recalling and post-thinking enable llms with long-term memory. *arXiv preprint arXiv:2311.08719*, 2023.

[62] Tomoyuki Kagaya, Thong Jing Yuan, Yuxuan Lou, Jayashree Karlekar, Sugiri Pranata, Akira Kinose, Koki Oguri, Felix Wick, and Yang You. Rap: Retrieval-augmented planning with contextual memory for multimodal llm agents. In *NeurIPS Workshop*, 2024.

[63] Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. Memorybank: Enhancing large language models with long-term memory. In *AAAI*, 2024.

[64] Hermann Ebbinghaus. Memory: A contribution to experimental psychology. *Annals of Neurosciences*, 2013.

[65] Gabriel Sarch, Yue Wu, Michael Tarr, and Katerina Fragkiadaki. Open-ended instructable embodied agents with memory-augmented large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, 2023.

[66] Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. Expel: Llm agents are experiential learners. In *AAAI*, 2024.

[67] Ali Modarressi, Ayyoob Imani, Mohsen Fayyaz, and Hinrich Schuetze. Ret-llm: Towards a general read-write memory for large language models. In *ICLR Workshop*, 2024.

[68] Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar prompting with memory for computer control. In *ICLR*, 2024.

[69] Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. A-mem: Agentic memory for llm agents. *arXiv preprint arXiv:2502.12110*, 2025.

[70] Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. A survey on the memory mechanism of large language model based agents. *arXiv preprint arXiv:2404.13501*, 2024.

[71] Yongchao Chen, Jacob Arkin, Yang Zhang, Nicholas Roy, and Chuchu Fan. Scalable multi-robot collaboration with large language models: Centralized or decentralized systems? In *ICRA*, 2024.

[72] Jonathan Light, Min Cai, Sheng Shen, and Ziniu Hu. From text to tactic: Evaluating llms playing the game of avalon. *arXiv preprint arXiv:2310.05036*, 2023.

[73] Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.

[74] Ge Zhang, Yemin Shi, Ruibo Liu, Ruibin Yuan, Yizhi Li, Siwei Dong, Yu Shu, Zhaoqun Li, Zekun Wang, Chenghua Lin, et al. Chinese open instruction generalist: A preliminary release. *arXiv preprint arXiv:2304.07987*, 2023.

[75] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

[76] Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, Shengyi Huang, Kashif Rasul, and Quentin Gallouédec. Trl: Transformer reinforcement learning. `https://github.com/huggingface/trl`, 2020.

[77] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *SOSP*, 2023.

[78] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.

# A    Supplementary Material

This supplementary material details the proposed method and presents additional experimental results. Section A.1 supplements the proposed method with preliminaries. Section A.2 presents more implementation details for experiments. Section A.3 reports more comparative results, and Section A.4 introduces extra ablation studies to validate the effectiveness of JoyAgents-R1. Section A.5 describes dataset configurations. Section A.6 includes all prompts used in baselines and our multi-agent architecture. Finally, Section A.7 provides extended case analyses. The code is included in *JoyAgents_R1_Code.zip*.

## A.1    Preliminaries

Multi-agent coordination faces great challenges due to the inherent complexity of the actor-critic framework, where the critic model with parameters comparable to the policy network exacerbates the computational overhead and convergence instability of typical MARL. To address these issues, this work adopts Group Relative Policy Optimization (GRPO), which foregoes the value function and computes advantages in a group-relative manner [24]. Given a query-answer pair $(q, a)$, the old policy $\pi_{\theta_{old}}$ of GRPO samples a group of outputs $\{o_i\}_{i=1}^G$. Then, the policy model $\pi_\theta$ is optimized by maximizing the objective as follows:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(O|q)]$$

$$\frac{1}{G} \sum_{i=1}^G \left( \min \left( \frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)} A_i, \text{clip} \left( \frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) \right) \tag{4}$$

where $\epsilon$ is a hyper-parameter for clipped objective like PPO [21]. Different from the original objective [23] which incorparates a KL penalty $\mathbb{D}_{KL}(\pi_\theta || \pi_{ref})$ between current and reference policies, our method removes this constraint to enable greater distributional divergence from the initial model following [73]. This operation saves memory usage during training by eliminating the need for multiple reference models, which is more suitable for the heterogeneous multi-agent coordination task. In addition, $A_i$ denotes the advantage of the $i$-th response by normalizing the group-wise rewards $\{r_i\}_{i=1}^G$ with the average and standard deviation as follows:

$$A_i = \frac{r_i - \text{mean}(\{r_1, r_2, \cdots, r_G\})}{\text{std}(\{r_1, r_2, \cdots, r_G\})} \tag{5}$$

## A.2    Implementation details

To improve the reproducibility of the experiments, more training settings for supervised fine-tuning (SFT) and reinforcement learning (RL) are provided in Table 3 and Table 4, respectively. To compare the SFT and RL methods fairly, the former is trained for a total of 10 epochs, while the latter is trained with SFT for 5 epochs followed by RL for 5 epochs. For RL training, the models are deployed on the TRL [76] framework on 2 NVIDIA H200 GPUs for accelerated inference via vLLM [77] and real-time weight updates, while allocating 6 GPUs for joint training.

## A.3    Comparative results

**Experimental setup for the Toolbench dataset.**    To validate the effectiveness of the proposed JoyAgents-R1, more comparative experiments are conducted on ToolBench [9]. The benchmark involves integrating API calls to accomplish tasks, where the agent should select the correct API and compose necessary API requests accurately. In this section, the test set is divided into in- and out-of-domain based on whether the tools used in the test instances have been seen during training. This setup enables us to evaluate both learning and generalization capabilities of the method. Moreover, the proposed method is compared with two baseline approaches, namely Single-think-SFT and Single-nothink-SFT. They are based on the Qwen2.5-3B model and fine-tuned on the toolbench training set. The former outputs the reasoning process and final results, while the latter directly generates the final response without the reasoning process. We also evaluate the performance of the larger open- and closed-source LLMs [1, 5, 23, 78] are also used to build agents for comparisons without fine-tuning.

Table 3: Training settings for supervised fine-tuning.

| Hyperparameter | Value |
|---|---|
| learning_rate | 5e-6 |
| max_length | 16384 |
| num_train_epochs | 5 |
| max_grad_norm | 1 |
| weight_decay | 0.01 |
| warmup_ratio | 0.03 |
| lr_scheduler_type | cosine |
| optim | adamw_torch |
| gradient_accumulation_steps | 4 |
| dataloader_num_workers | 8 |
| per_device_train_batch_size | 1 |

Table 4: Training settings for reinforcement learning.

| Hyperparameter | Value |
|---|---|
| learning_rate | 1e-6 |
| max_length | 16384 |
| num_train_epochs | 5 |
| max_grad_norm | 1 |
| adam_epsilon | 1e-5 |
| num_groups | 5 |
| topk_groups | 5 |
| kl_coef | 0 |
| grpo_epoch | 2 |
| policy_clip_eps | 0.2 |
| temperature | 1.2 |
| per_device_train_batch_size | 1 |

**More metrics on the Toolbench dataset.** To comprehensively compare the model performance, a variety of evaluation indicators are used as follows:

- Plan_ACC: The accuracy of the agent's planning decisions at each step of the tool calling.

- Act_EM: The proportion of predicted API names that exactly match the real API names.

- Easy_F1: The predicted argument F1 score when the ground truth argument is empty.

- Hard_F1: The predicted argument F1 score when the ground truth argument is not empty.

- F1: The predicted argument F1 score across all conditions.

- No_Hallu: The frequency of predicted API names that do not have hallucinations.

- Avg: The average value of the above indicators.

**More comparative results on the Toolbench dataset.** As illustrated in Table 5, since the comparison results under the in-domain and out-of-domain settings show roughly the same trend, they are analyzed from two perspectives as follows:

- *Compare with closed-source and open-source models with larger parameters.* Compared to DeepSeek-R1 and DeepSeek-V3, JoyAgents-R1 exceeds them in terms of Act_EM, Hard_F1, F1, and Average. For instance, DeepSeek-R1 is relatively improved by 7.3%, 73.1%, 36.8%, and 8.7% in the above indicators under in-domain scenarios, while 10.6%, 81.4%, 37.4%, and 11.0% in the above indicators under out-of-domain scenarios, demonstrating its superiority in function calling. It is worth mentioning that the number of parameters activated during the inference of DeepSeek-R1/V3 is over 37B, while our model achieves a more satisfying performance with 3B parameters. Meanwhile, compared to larger Qwen2.5 series models (*e.g.*, 14B and 32B parameters), our 3B model significantly outperforms by using multi-agent training frameworks, demonstrating

15

Table 5: More comparative results (%) on the Toolbench dataset. Bold represents the optimal score, and '_' represents the suboptimal value. The higher the value, the better the performance.

| Method | Plan_ACC | Act_EM | Easy_F1 | Hard_F1 | F1 | No_Hallu | Avg |
|---|---|---|---|---|---|---|---|
| **In the domain** | | | | | | | |
| GPT-4o | **82.5** | **57.0** | **31.1** | 22.2 | _37.3_ | _99.9_ | **55.0** |
| DeepSeek-R1 | 72.3 | 45.2 | _18.7_ | 18.2 | 28.5 | **100.0** | 47.1 |
| DeepSeek-V3 | 68.4 | 38.0 | 14.9 | 14.3 | 23.6 | 99.4 | 43.1 |
| Qwen2.5-32B | 72.7 | 43.7 | 15.1 | 21.0 | 28.8 | 99.3 | 46.8 |
| Qwen2.5-14B | 42.2 | 15.0 | 3.7 | 7.1 | 10.1 | **100.0** | 29.7 |
| Single-nothink-SFT | 55.7 | 33.0 | 10.1 | 20.5 | 26.5 | 95.7 | 40.3 |
| Single-think-SFT | 73.3 | 47.0 | 16.3 | _29.6_ | 36.8 | 94.5 | 49.6 |
| JoyAgents-R1 (3B) | _73.5_ | _48.5_ | 17.6 | **31.5** | **39.0** | 96.7 | _51.2_ |
| **Out of the domain** | | | | | | | |
| GPT-4o | **83.9** | **63.2** | **42.9** | 22.3 | 44.1 | 99.8 | **59.3** |
| DeepSeek-R1 | 74.6 | 52.0 | 28.3 | 20.4 | 35.6 | 99.8 | 51.8 |
| DeepSeek-V3 | _77.2_ | 56.1 | 32.0 | 21.9 | 38.2 | **100.0** | 54.2 |
| Qwen2.5-32B | 28.6 | 0.0 | 0.0 | 0.0 | 0.0 | **100.0** | 21.4 |
| Qwen2.5-14B | 43.9 | 16.6 | 7.2 | 6.2 | 12.3 | _99.9_ | 31.0 |
| Single-nothink-SFT | 59.7 | 41.1 | 19.7 | 25.9 | 36.5 | 97.3 | 46.7 |
| Single-think-SFT | 72.3 | 57.0 | _32.3_ | _36.1_ | _48.7_ | 96.3 | 57.1 |
| JoyAgents-R1 (3B) | 73.4 | _57.5_ | 31.1 | **37.0** | **48.9** | 96.9 | _57.5_ |

Table 6: Ablation results (%) in updating top-$K$ models. 'FC' is the function call.

| Method | Math | QA | E-commerce FC | General FC | Cooperation | Average |
|---|---|---|---|---|---|---|
| Top-1 | 67.0 | 16.0 | **48.0** | 75.0 | 4.0 | 42.0 |
| Top-2 | 57.0 | 12.0 | 41.0 | 75.0 | **6.0** | 38.2 |
| Top-All | 64.0 | 17.0 | 44.0 | 72.0 | 3.0 | 40.0 |
| Top-5 | **68.0** | **22.0** | **48.0** | **76.0** | **6.0** | **44.0** |

that even smaller language models can achieve competitive results through our joint evolution approach. In addition, it can be observed that JoyAgents-R1 achieves the best results in Hard_F1 and F1 under both in-domain (*i.e.*, 31.5% and 39.0%) and out-of-domain settings (*i.e.*, 37.0% and 48.9%), even surpassing GPT-4o, showing the great potential of applying RL to improve small open-source models.

- *Compare with baseline methods.* JoyAgents-R1 demonstrates superior performance over baseline methods, particularly in Plan ACC, Act_EM, F1, and Average metrics. For example, the average score of Single-nothink-SFT under in-domain is increased from 40.3% to 49.6% by enabling the think process. Then, the average value of Single-think-SFT is further improved by 3.2% through multi-agent reinforcement learning. This validates two key advantages:

  – The framework's strength in advancing agent planning and API calling capabilities.
  – The effectiveness of the joint evolution dynamics training strategy in refining individual agent decision-making through coordinated policy adaptation.

## A.4 Ablation study

**Ablation on updating top-$K$ models.** Table 6 presents the performance outcomes of updating varying numbers of nodes along a trajectory. The empirical results reveal that updating the top-5 nodes yields optimal performance, outperforming alternative strategies. Specifically, compared to updating only the top-1 node, the top-5 update achieves a 37.5% improvement on the QA dataset and a 50% gain in collaborative tasks. Relative to updating top-2 nodes, the top-5 approach delivers a 19.3% boost on the math dataset and an 83.3% enhancement on QA tasks. compared to updating all nodes, updating the top-5 nodes results in a 29.4% improvement on QA and a 100% increase in

Table 7: Ablation results (%) on the number of sub-agents. 'FC' is the function call.

| Method | QA | E-commerce FC | Math | Cooperation |
|---|---|---|---|---|
| Master agent + 2 sub-agents | **26.0** | **55.0** | - | - |
| Master agent + 3 sub-agents | 25.0 | 51.0 | **69.0** | 5.0 |
| Master agent + 4 sub-agents | 22.0 | 48.0 | 68.0 | **6.0** |

collaborative task success rates. These findings demonstrate that merely updating the top-1 or top-2 nodes is insufficient for holistic system optimization, as such localized adjustments fail to address systemic weaknesses. Conversely, updating all nodes lacks global awareness and focus, leading to redundant computations. Our top-5 strategy, however, employs global trajectory analysis to identify and update the weakest nodes, *i.e.*, those most limiting system performance, thereby maximizing efficiency and efficacy. This selective updating mechanism ensures that optimization efforts are concentrated on critical bottlenecks, yielding superior overall performance.

**Ablation on the number of sub-agents.** Table 7 presents the performance of the master agent when integrated with varying numbers of sub-agents across different datasets. Specifically, the configurations include:

- *2 sub-agents*: QA agent and e-commerce function-call agent.
- *3 sub-agents*: The aforementioned two plus the math agent, with the addition of collaborative tasks.

The results indicate that for individual tasks, configurations with fewer sub-agents (2 or 3) outperform the full set of sub-agents. For instance, on the QA dataset, the 2-sub-agent setup yields an 18.2% improvement, while the 3-sub-agent setup achieves a 13.6% gain. Similarly, on the e-commerce function-call dataset, the 2-sub-agent and 3-sub-agent configurations exhibit 14.6% and 6.2% improvements, respectively. These findings align with our intuition that fewer sub-agents are more effective for non-collaborative tasks, due to reduced complexity and interference. Conversely, in collaborative tasks, the full set of 4 sub-agents demonstrates superior performance, attributed to their exposure to a broader range of data and interactions. This outcome underscores the efficacy of our multi-agent system design, which is specifically tailored to address collaborative challenges. The enhanced performance in collaborative scenarios validates the structural design of our system, highlighting the benefits of a comprehensive multi-agent framework in handling complicated and interdependent tasks.

## A.5 Datasets and setup

The datasets used in this paper can be divided into two categories: SFT and RL. Since the RL stage is trained end-to-end, this type of data only contains the initial query and the final response. More details of the SFT dataset used for each agent are introduced as follows:

**Master agent datasets.** As shown in Table 8, the case for the master agent includes user queries, optional tools, invoked agents, retrieved memories, historical dialogues, and tool-generated responses. In addition, the reasoning processes (*i.e.*, $< think > \cdots < /think >$), tool calling (*i.e.*, $< tool\_call > \cdots < /tool\_call >$), and final answers are provided in the dataset.

**Function-call agent datasets.** Regarding function-calling, the e-commerce agent incorporates 12 prevalent APIs of e-commerce platforms, while the general agent draws from 1,000 diverse function calls from ToolBench [9]. Table 9 illustrates a case for the e-commerce function-call agent, which is called the 'rody agent' in our multi-agent system for easy reference in code and expressions.

**Question-answering agent datasets.** The question-answering agent includes 'tool_retrievals' and 'tool_summary' as demonstrated in Table 10. Its data covers open-domain [74] and e-commerce scenarios. In this work, the QA agent has an alias called 'expert agent' to show its professionalism.

**Math agent datasets.** The dataset for the math agent is constructed from GSM8K [75]. The last part of Table 8 shows the thinking process and answers of math agents.

Table 8: A case of the master agent in a collaborative task.

| | Case #1 of the master agent in a collaborative task |
|---|---|

**User:** I want to open a small personal store on [e-commerce platform A] to sell educational toys. How much deposit do I need to pay? If I open two stores, how much deposit will I need to pay in total?

| | |
|---|---|
| **System** | **{Role}** from Template #1 in Section A.6<br>**Specialized Retrieval Tools/Memory:**<br>< ***tools*** ><br>[{"name": "math_agent", "description": "When the user puts forward a clear request to give a solution or answer to a specific mathematical problem, the agent needs to be called. If it is only about a mathematically related topic, the agent cannot be called", "parameters": {}}, ...]<br>< ***/tools*** ><br>< ***memory*** ><br>[{"ID": 1620, "Query": " Calculate how many times Danny folded his picnic blankets given the initial and final area.", "Plan": ["math_agent"], "Answer": "How many square feet are each blanket? ** The total area of each blanket is 8*8=«8*8=64»64. How many square feet are the blankets? ** Their total area is 192 square feet because 3 x 64 = «3*64=192»192. How many times did he fold them? ** He folded them 4 times because 192 / 48 = «192/48=4»4. #### 4", "Time": "2025-03-26 11:41:13", "Score": 6} ...]<br>< ***/memory*** ><br>**{Problem Resolution Framework}** from Template #1 in Section A.6 |
| **Assistant** | < ***tool_call*** ><br>{"name": "rody_agent", "arguments": {"intention": ["Query the deposit required for opening a small personal store on [e-commerce platform A] to sell educational toys."]}}<br>< ***/tool_call*** > |
| **Tool** | Pass calling {"api_name": "check_shop_expenses", "parameter": {"keyword": ["Educational Toys"], "shop_body": ["Individual"], "shop_type": ["small shops"]}}, the API call was successful.<br>Results are as follows:<br>1. Specific amount of margin per year: 6000 yuan.<br>2. Transaction service fee rate per transaction: 0.25%.<br>3. Technical service fee rate per transaction: 0.41%. |
| **Assistant** | < ***tool_call*** ><br>{"name": "tool_retrievals", "arguments": {"intention": ["Calculate the total deposit required for opening two small personal stores on [e-commerce platform A] to sell educational toys, based on the queried deposit information of 6000 yuan per store."]}}<br>< ***/tool_call*** > |

| | |
|---|---|
| **Tool** | [{"name": "expert_agent", "description": "When users chat, greet, abuse, make jokes, discuss music/games/current politics/military/entertainment/food and other daily topics, or consult the rules of the e-commerce platform, the agent needs to be invoked ", "parameters": {}}, {"name": "math_agent", "description": "When the user puts forward a clear request to give a solution or answer to a specific mathematical problem, the agent needs to be called. If it is only about a mathematically related topic, the agent cannot be called ", "parameters": {}}, {"name": "toolbench_agent", "description": "function-0: This is the subfunction for tool "football_xg_statistics" and its function is "seasons list by tournamentid for football xg statistics", you can use this tool. The description of this function is: Football (soccer) xG statistics. We provide statistics for more than 80 leagues. There are more than 90,000 games and 2 million shots in our database. The method is used to get a list of seasons by the league ID; function-1: This is the subfunction for the tool "mailslurp_email_testing" and its function is "getbouncedrecipients for mailslurp email testing", you can use this tool. The description of this function is: Email sandbox testing API. Bounced recipients are email addresses that you have sent emails to that did not accept the email. Once a recipient is bounced, you cannot send emails to that address."}] |
| **Assistant** | $< tool\_call >$<br>{"name": "math_agent", "arguments": {"intention": ["Calculate the total deposit required for opening two small personal stores on [e-commerce platform A] to sell educational toys, based on the queried deposit information of 6000 yuan per store."]}}<br>$< /tool\_call >$ |
| **Tool** | The user wants to open a small personal store on [e-commerce platform A] to sell educational toys. The deposit required for one store is 6000 yuan per year. If the user opens two stores, the total deposit will be 6000 yuan multiplied by 2. 12000 |
| **Assistant Final** | $< think >$<br>The user's issue is to calculate the total deposit required for opening two small personal stores on [e-commerce platform A] to sell educational toys. The tool rody_agent was initially called to query the deposit information, which returned that the deposit for one store is 6000 yuan per year. The user then asked for the total deposit for two stores. The tool math_agent was called to perform the calculation, which returned the result of 6000 yuan multiplied by 2, resulting in 12000 yuan. Since the tool math_agent successfully calculated the total deposit, the final solution can be directly provided based on the tool's response. Therefore, the final solution is "12000".<br>$< /think >$<br>12000 |

Table 9: A case of the function-call agent in an e-commerce task.

| | **Case #2 of the function-call agent in an e-commerce function-call task** |
|---|---|

**User:** Hello, I found that some orders may have the label of door-to-door replacement, what information should be provided in this case?

| | |
|---|---|
| **System** | **{Role}** from Template #2 in Section A.6 <br> **Specialized Retrieval Tools/Memory:** <br> $< tools >$ <br> [{"name": "search_payment_method","description": "Inquire the settlement method of the user's shop, such as B2B settlement, C2B settlement, bank card settlement, online banking wallet settlement.","parameters": {}}, {"name": "search_order_code", "description": "...", ,"parameters": {...}}, ...] <br> $< /tools >$ <br> $< memory >$ <br> [{"ID": "1254", "Query": "Hello, I found that some orders may have the label of door-to-door replacement, what information should be provided in this case?", "Plan": [], "Answer": "Pass calling {'api_name': 'search_order_code', 'parameter': {'order_service_mix_tag': ['door-to-door replacement'], 'search_keys': ['order label']}}, the API call was successful. <br> Results are as follows: <br> We have retrieved the relevant order information for you.", "Time": "2025-04-29 02:20:30", "Score": 10}, ...] <br> $< /memory >$ <br> **{Problem Resolution Framework}** from Template #2 in Section A.6 |
| **Assistant** | $< tool\_call >$ <br> {"name":"search_order_code", "arguments":{"order_service_mix_tag": ["door-to-door replacement"], "search_keys": ["order label"]}} <br> $< /tool\_call >$ |
| **Tool** | Pass calling {"api_name": "search_order_code", "parameter":{"order_service_mix_tag": ["door-to-door replacement"], "search_keys": ["order label"]}}, the API call was successful. <br> Results are as follows: <br> We have retrieved the relevant order information for you. |
| **Assistant Final** | $< think >$ <br> The Input issue is about retrieving information for orders with the label "door-to-door replacement." The Output is the result of calling the "search_order_code" tool with the parameters '{"order_service_mix_tag": ["door-to-door replacement"], "search_keys": ["order label"]}'. The Output indicates that the API call was successful and relevant order information has been retrieved. This matches the reference information in <memory>, where a similar query was resolved by calling the same tool with the same parameters. Therefore, the Output is directly derived from the reference information, and the reasoning process is consistent with the provided solution. <br> $< /think >$ <br> Pass calling {"api_name": "search_order_code", "parameter": {"order_service_mix_tag": ["door-to-door replacement"], "search_keys": ["order label"]}}, the API call was successful. <br> Results are as follows: We have retrieved the relevant order information for you. |

Table 10: A case of the question-answering agent in a QA task.

---

**Case #3 of the question-answering agent in a QA task**

---

**User:** Where to open the LCL service?

---

| | |
|---|---|
| **System** | {**Role**} from Template #4 in Section A.6<br>**Specialized Retrieval Tools/Memory:**<br>$< tools >$<br>[{"name": "tool_retrievals", "description": "Knowledge retrieval tool. Searches for potentially relevant information in a vertical knowledge base (e-commerce merchant operations scenarios) based on the current question.", "arguments": {}} {"name": "tool_summary", "description": "Intelligent Q&A tool. When retrieval content is provided, it primarily answers based on the retrieved content; otherwise, it responds based on its knowledge.", "arguments": {}}]<br>$< /tools >$<br>$< memory >$<br>[{"id": 4439, "model_name": "expert", "query": "Is my shop bounded with a card?", "plan": [], "response": "You can check and manage your bank card information in the backstage of your store, in the "store management"-"store information"-"settlement configuration".", "score": 1.0, "time": 7514}, ...]<br>$< /memory >$<br>{**Problem Resolution Framework**} from Template #4 in Section A.6 |
| **Assistant** | $< tool\_call >$<br>{"name": "tool_retrievals", "arguments": {}}<br>$< /tool\_call >$ |
| **Tool** | [General] Corporation_A Merchant Consolidation Warehouse<br>Opening Configuration Page (see link: ...)<br>[Other] Merchant Question: Consolidation Opening Entry<br>Merchant Assistant:<br>1) How to open: You can click on the left menu bar of the Merchant Center - Services - Merchant Service Center - Sales - Remote Area Logistics Service - Click to open to enter the product configuration page. Follow the page prompts to configure your products, and then turn on the service switch in the upper right corner to take effect. Direct link: ...<br>2) Fees: For products bound to consolidation and free shipping, you only need to bear the cost of the first leg to Xi'an. The second leg from Xi'an to the consumer is fully subsidized by the platform, and consumers can enjoy free shipping and free door-to-door delivery.<br>3) Benefits: Free second-leg logistics for remote area shipments, free interception for second-leg logistics in remote areas, estimated doubling of remote area order volume, and exemption from store comprehensive experience scores. ... |
| **Assistant** | $< tool\_call >$<br>{"name": "tool_summary", "arguments": {}}<br>$< /tool\_call >$ |
| **Tool** | The entry to open the consolidation service is as follows:<br>1. Click the left menu bar.<br>2. Select "Service".<br>3. Enter "Merchant Service Center".<br>4. Select "Remote Area Logistics Service" under "Selling Section".<br>5. Click "Open" to enter the product configuration page.<br>6. Configure the products according to the page prompts, and turn on the service switch in the upper right corner to take effect.<br>Direct link: ... |

| | |
|---|---|
| **Assistant Final** | < *think* > <br> Given that the user's question pertains to e-commerce-specific knowledge and both tool_retrievals and tool_summary have already been invoked in the historical context, the optimal path is to directly output the answer provided by tool_summary. This is because the necessary information has been retrieved and summarized, and there's no need for further tool invocations, ensuring efficiency and cost-effectiveness in tool usage. <br> < */think* > <br> The entry to open the consolidation service is as follows: <br> 1. Click the left menu bar. <br> 2. Select "Service". <br> 3. Enter "Merchant Service Center". <br> 4. Select "Remote Area Logistics Service" under "Selling Section". <br> 5. Click "Open" to enter the product configuration page. <br> 6. Configure the products according to the page prompts, and turn on the service switch in the upper right corner to take effect. <br> Direct link: ... |

## A.6 Prompts

**System prompts for multi-agents.** The system prompt template for each agent in our multi-agent architecture is as follows:

---

**Template #1: System Prompt for the Master Agent**

**Role:**

You are a customer service expert of an e-commerce platform, specializing in selecting appropriate tools and agents to solve user problems based on user questions. Please understand and analyze the user's current problem according to the history dialogue until the user's problem is solved. There are some tools available between < *tools* >< */tools* >for selection at each step. You can think as follows:

- If the user's intention can completely match other known specific tools except for tool_intention and tool_retrievals, directly call that tool.

- If you cannot understand the user's current intention, call tool_intention.

- If the user's intention is clear but cannot match other known specific tools except for tool_intention and tool_retrievals, call tool_retrievals to get complementary information to choose from all tools again.

- If an agent tool has been called and returned results, or content has been recalled between < *memory* >< */memory* >, you need to judge whether the current agent result or the similar answers recorded in the memory meet the user's intention. If so, directly output the answer based on it; otherwise, re-arrange.

**Specialized Retrieval Tools / Memory:**

< *tools* >

- {"name": "rody_agent", "description": "The agent has the following functions:

  - Check fee: check the deposit, quality guarantee, and technical service fee requirements required for settling in the Corporation_A platform;

  - Check qualifications: inquire about the various documents/material requirements required for settling in the Corporation_A platform;

  - Check the order: Check all questions about the related order;

  - Query after-sales information: Check all questions about the related after-sales order;

---

- Order reporting: Delayed reporting of orders that cannot be delivered on time, but unable to check whether the order has been reported or whether it needs to be reported;
- Check the refund of deposit at the time of check-out: When the user applies for check-out and returns the deposit, check the current refund progress of the deposit and the reason why the deposit cannot be returned at present;
- Query product promotion: Quickly obtain detailed information related to the designated promotion activities of the merchants;
- Query product coupons: obtain the status of coupons specified by the merchant and applicable products, but can not query the reason why the coupons are not effective;
- Check the status of bank card: Help users check the current binding and verification status of their bank card;
- Query the settlement method: According to the actual situation of the user, help query the settlement method of its store, and provide the current store payment and platform refund flow;
- Query product audit: According to the actual requirements of the user, query the review status of the user's product listing/modification and the reasons for the slow review progress, but can not query the reasons for the failure of the product audit."}

- {"name": "expert_agent", "description": "When users chat, greet, abuse, make jokes, discuss music/games/current politics/entertainment/food and other daily topics, or consult the rules of the e-commerce platform, the agent needs to be invoked.", "parameters": {}}

- {"name": "math_agent", "description": "When the user puts forward a clear request to give a answer to a specific mathematical problem, the agent needs to be called. If it is only about a mathematically related topic, the agent cannot be called.", "parameters": {}}

- {"name": "toolbench_agent", "description": "", "parameters": {}}

- {"name": "tool_intention", "description": "The intention understanding tool. Understand the user's real intention based on context and current question.", "arguments": {}}

- {"name": "tool_retrievals", "description": "The API retrieval tool. Retrieve related APIs from the API knowledge base based on intention.", "arguments": {"intention": "user's current intention"}}

- ...

$< /tools >$
$< memory >$
memory_append(Optional)
$< /memory >$

**Problem Resolution Framework:**

Specific requirements are as follows:

- When selecting tools, please refer to the tool's function description. Each tool's function only contains the content in the description, and it is prohibited to guess or extend other functions based on the description.

- When selecting tools after calling tool_retrievals, you can only choose from the candidate tool set retrieved through tool_retrievals. Each tool's function only contains the content in the description.

- Please be faithful to the semantics of the current problem and historical dialogue, and do not output content that does not exist in the historical dialogue and current problem.

- Directly output the agent's response if there is no error.

- Do not alter or truncate any words if the response is from rody/math/toolbench.

- For math or coding questions, use the user's query as the sub-agent calling intention.

- Do not call the same tool repeatedly. If the result from the previous tool call is incorrect, try using other tools.

- When a user requests a solution or answer to a specific math problem, output in the following format (output the numerical answer directly after $</think>$, without any units or irrelevant characters): $<think>$think process$</think>$answer.
- Output strictly according to the following format:
  - The user's question is not clear, or unable to understand the user's intention:
    * $<think>$ The analysis and thinking process of the user's problem, and the reason for calling the intention recognition tool. $</think>$
    * $<tool\_call>$ {"name": "tool_intention", "arguments": {}} $</tool\_call>$
  - The tools in <tools></tools> and all the currently retrieved tools cannot meet the user's intention:
    * $<think>$ The analysis and thinking process of the user's problem, and the reason for calling the tool_retrievals tool. $</think>$
    * $<tool\_call>$ {"name": "tool_retrievals", "arguments": {"intention": user_intention}} $</tool\_call>$
  - The agent tools in <tools></tools> or all the currently retrieved tools can meet the user's intention:
    * $<think>$ The analysis and thinking process of the user's problem, and the reason for choosing which agent to solve the user's problem. $</think>$
    * $<tool\_call>$ {"name": chosen agent name, "arguments": {"intention": user's intention that can be met through chosen agent}}$</tool\_call>$
  - According to the recalled information between $<memory></memory>$ or the agent calling result, you can answer the user's question:
    * $<think>$ The reason why you can answer the user's problem based on the current known information.$</think>$
    * The answer to the current user's question.

---

## Template #2: System Prompt for the E-commerce Function-call Agent

**Role:**

You are a customer service expert of Corporation_A e-commerce platform, specializing in solving user problems based on user questions and selecting the final API tools. Please understand and analyze the user's current problem according to the history dialogue, thinking step by step until the user's problem is solved. There are some tools available for selection at each step. You can think as follows:

- If the user's intention can completely match other known specific APIs except for tool_intention and tool_retrievals, directly call that API and identify the required parameters.
- If the user's intention is clear but cannot match other known specific APIs except for tool_intention and tool_retrievals, call tool_retrievals to recall some related APIs to choose again.
- If an API has been called and returned results, or content has been recalled between $<memory></memory>$, you need to judge whether the current API result or the similar answers recorded in the memory meet the user's intention. If so, directly output the answer based on it; otherwise, re-arrange.

**Specialized Retrieval Tools / Memory:**

$<tools>$

- {"name": "tool_retrievals", "description": "Knowledge retrieval tool. Searches for potentially relevant information in a vertical knowledge base (e-commerce merchant operations scenarios) based on the current question.", "arguments": {}}

- ...

$</tools>$
$<memory>$
memory_append(Optional)
$</memory>$

**Problem Resolution Framework:**

Specific requirements are as follows:

- When selecting tools, please refer to the tool's function description. Each tool's function only contains the content in the description, and it is prohibited to guess or extend other functions based on the description.

- When selecting APIs, you can only choose from the candidate API set recalled through tool_retrievals. Each API's function only contains the content in the description.

- Please be faithful to the semantics of the current problem and historical dialogue, and do not output content that does not exist in the historical dialogue and current problem.

- Do not call the same tool more than once, and try to call different APIs.

- Just output the tool response directly if there is no error.

- Output strictly according to the following format:

  - The APIs in $<tools></tools>$ and all the currently recalled APIs cannot meet the user's intention: $<think>$ The analysis and thinking process of the user's problem, and the reason for calling the API retrieval tool. $</think><tool\_call>$ {"name": "tool_retrievals", "arguments": {"intention": user_intention}} $</tool\_call>$

  - The APIs in $<tools></tools>$ or all the currently recalled APIs can meet the user's intention: $<think>$ The analysis and thinking process of the user's problem, and the reason for choosing which api to solve the user's problem. $</think>$ $<tool\_call>$ {"name": chosen api name, "arguments": parameters passed to api} $</tool\_call>$

  - According to the recalled information between $<memory></memory>$ or the API calling result, you can answer the user's question:

    * $<think>$ The reason why you can answer the user's problem based on the current known information.$</think>$
    * The answer to the current user's question.

---

### Template #3: System Prompt for the General Function-call Agent

**Role:**

You are a customer service expert of Corporation_A e-commerce platform, specializing in solving user problems based on user questions and selecting the final API tools. Please understand and analyze the user's current problem according to the history dialogue, thinking step by step until the user's problem is solved. There are some tools available for selection at each step. You can think as follows:

- If the user's intention can completely match other known specific APIs except for tool_intention and tool_retrievals, directly call that API and identify the required parameters.

- If the user's intention is clear but cannot match other known specific APIs except for tool_intention and tool_retrievals, call tool_retrievals to recall some related APIs to choose again.

- If an API has been called and returned results, or content has been recalled between $<memory></memory>$, you need to judge whether the current API result or the

similar answers recorded in the memory meet the user's intention. If so, directly output the answer based on it; otherwise, re-arrange.

**Specialized Retrieval Tools / Memory:**

$< tools >$

- {"name": "tool_retrievals", "description": "The API retrieval tool. Retrieve related APIs from the API knowledge base based on intention.", "arguments": {"intention": "user's current intention"}}

- ...

$< /tools >$
$< memory >$
memory_append(Optional)
$< /memory >$

**Problem Resolution Framework:**

Specific requirements are as follows:

- When selecting tools, please refer to the tool's function description. Each tool's function only contains the content in the description, and it is prohibited to guess or extend other functions based on the description.
- When selecting APIs, you can only choose from the candidate API set recalled through tool_retrievals. Each API's function only contains the content in the description.
- Please be faithful to the semantics of the current problem and historical dialogue, and do not output content that does not exist in the historical dialogue and current problem.
- Do not call the same api more than once, and try to call different APIs.
- Just output the tool response directly if there is no error, or there is no other appropriate api to call.
- Output strictly according to the following format:
  - The APIs in $< tools >< /tools >$ and all the currently recalled APIs cannot meet the user's intention: $< think >$ The analysis and thinking process of the user's problem, and the reason for calling the API retrieval tool. $< /think >< tool\_call >$ {"name": "tool_retrievals", "arguments": {"intention": user_intention}} $< /tool\_call >$
  - The APIs in $< tools >< /tools >$ or all the currently recalled APIs can meet the user's intention: $< think >$ The analysis and thinking process of the user's problem, and the reason for choosing which api to solve the user's problem. $< /think >$ $< tool\_call >$ {"name": chosen api name, "arguments": parameters passed to api} $< /tool\_call >$
  - According to the recalled information between $< memory >< /memory >$ or the API calling result, you can answer the user's question:
    * $< think >$ The reason why you can answer the user's problem based on the current known information.$< /think >$
    * The answer to the current user's question.

---

Template #4: System Prompt for the QA Agent

**Role:**

You are a customer service expert for an e-commerce platform, capable of utilizing your memory and searching for appropriate tools to address user inquiries. Based on the current

question and past tool selections and their responses, proceed step-by-step to determine which tool to use or what content to output next.

**Specialized Retrieval Tools / Memory:**

1. Tools at your disposal: Results from the same tool & arguments are unique.
$< tools >$
- {"name": "tool_retrievals", "description": "Knowledge retrieval tool. Searches for potentially relevant information in a vertical knowledge base (e-commerce merchant operations scenarios) based on the current question.", "arguments": {}}
- {"name": "tool_summary", "description": "Intelligent Q&A tool. When retrieval content is provided, it primarily answers based on the retrieved content; otherwise, it responds based on its knowledge.", "arguments": {}}

$< /tools >$
2. Your memory content: Memory varies for different questions. In memory, the shorter the plan route and the higher the score, the more valuable it is for reference.
$< memory >$
memory_append(Optional)
$< /memory >$

**Problem Resolution Framework:**

- Requirements for Tool Selection:
  - tool_retrievals: Call this tool when retrieving domain-specific knowledge related to merchants or e-commerce scenarios is required.
  - tool_summary: Call this tool when generating a response to the user's question is needed.
    * tool_summary will refer to the results from tool_retrievals to generate a response only when the previous call was to tool_retrievals; otherwise, it will respond directly.
  - Outputting Answers: (Must Pay Attention!) You cannot answer questions directly. The following scenarios apply when outputting answers:
    * If a similar question exists in the memory, you can directly output the answer provided in the memory without invoking any tools.
    * If no usable answer is found in the memory, you must first call tool_summary and return its output result as is (without modifying the output of tool_summary).
  - Tool Efficiency: Tool invocations incur costs. If the required information is sufficient to answer the user's question, respond directly without unnecessary tool calls.

- Special Handling: When the question contains content such as "Pass calling ... Results are as follows ...", this part represents the results of historical API calls. You do not need to answer this part in your response. However, when providing the final answer, you must combine the tool_summary response with the historical API call results and include them together. In other words, after using the tool_summary to answer the question, add the historical API call results to the beginning of the tool_summary response and return them together.

- Applying the Above Tool Selection Requirements: When selecting tools for the current round, consider the following:
  - Analyze Previously Called Tools:
    * Important: Avoid calling the same tool that has already been called in history.
  - (If memory is not empty) Analyze memory content:
    * If the user's question is essentially identical to one in memory, output the answer from memory without calling other tools.
    * If the user's question is similar in content or type to one in memory, refer to the plan in memory for guidance. For example, if the user's question and the memory question both pertain to vertical domain knowledge, you can follow the plan in memory.

> ∗ If the user's question bears no similarity to the memory content, ignore the memory.
> – (After confirming memory does not provide a direct answer) Determine whether the user's question requires e-commerce or Corporation_A-specific knowledge to answer:
>   ∗ If no vertical domain knowledge is needed, call tool_summary directly. Otherwise, first call the tool_retrievals to retrieve relevant knowledge before answering.
> – If the Previous Call Was to tool_summary:
>   ∗ Since tool_summary cannot be called again, and the output must come from tool_summary, directly output the answer.
>
> • Output format requirements:
>
> – Answer output format: $<think>$ Thought process $</think>$ Output answer(from tool_summary or memory; Important: Any API call content found in the question must be included verbatim in the final response.)
> – Tool call format: $<think>$ Thought process $</think>$ $<tool\_call>$ {"name": "tool_name", "arguments": {"param": "value"}}$</tool\_call>$

---

**Template #5: System Prompt for the Math Agent**

**Role:**

You are a math expert, specializing in step-by-step thinking to answer the math problems raised by users. Now you have a memory library, and the relevant memories will be stored in it. You can combine the content in the memory to answer questions. The specific thinking steps are as follows:

• If there is an identical question in the memory, you can use the answer of that question to directly answer the current question.

• If all the questions in the memory are different from the current user's question, you need to think and answer by yourself.

**Specialized Retrieval Tools / Memory:**

$<tools>$
[tool_append](Optional)
$</tools>$
$<memory>$
memory_append(Optional)
$</memory>$

**Problem Resolution Framework:**

According to different situations, the output should strictly follow the following format:

• If there is an identical question in the memory:

– $<think>$ The reason for choosing the answer to the identical question.$</think>$
– The answer to the current question.

• If there is no identical question in the memory:

– $<think>$ The reason for not choosing a question from memory, and the steps of thinking about the current user's question.$</think>$
– The answer to the current question.

**System prompt for the single agent.** The following is a system prompt template for single-agent multi-step reasoning based on open-source or closed-source SOTA models:

## Template #6: System Prompt for the Single Agent

**Role:**

You are a customer service expert of an e-commerce platform(Corporation_A), specializing in selecting appropriate tools and agents to solve user problems based on user questions. Please understand and analyze the user's current problem according to the historical information until the user's problem is solved. There are some tools available between <tools></tools>for selection at each step.

**Specialized Retrieval Tools:**

$< tools >$

- {"name": "tool_retrievals_knowledge", "description": "Vertical knowledge base search tool (e-commerce merchant operations context). Identifies relevant information based on user queries.", "arguments": {"intention": "user's current intention or query"}}

- {"name": "tool_retrievals_API_shop", "description": "E-commerce platform API lookup. Retrieves relevant APIs from the API knowledge base using intent analysis.", "arguments": {"intention": "user's current intention or query"}}

- {"name": "tool_retrievals_API_general", "description": "General API lookup. Retrieves relevant APIs from the API knowledge base using intent analysis.", "arguments": {"intention": "user's current intention or query"}}

$< /tools >$

**Problem Resolution Framework:**

1. Question Types & Response Protocols:
You may encounter different types of questions. The types of questions and the required output formats are as shown below:

- Math problems:
  - Provide direct solutions to numerical queries.
  - Output in the following format (Provide the numerical answer directly after $< /think >$, without units or any irrelevant characters): $< think >...< /think >$Final numeric answer

- API scheduling problems:
  - The APIs are divided into e-commerce platform APIs and general APIs.
  - When API tools are required: Use relevant tool_retrievals to identify candidate APIs (original/paraphrased queries accepted).
  - Output the API call results in the following format: $< tool\_call >$"name": "API_name", "arguments": "key1":["value11", "value12"], "key2":["value21", "value22"]...$< /tool\_call >$
  - Some solutions require sequential API calls, but you can just call only one API at each step. Use prior outputs as inputs for subsequent calls.

- Q&A problems:
  - Engage directly in casual conversations (greetings/jokes/daily topics).
  - For e-commerce policy queries: Invoke tool_retrievals_knowledge for domain knowledge. Respond based on retrieved content.

2. Tool/API Selection Guidelines:

- The results of the previous Tool/API call will be returned in the format $< tool\_response >...< /tool\_response >$.

- The response format for API dispatching results is: "Pass calling ... Results are as follows: ...". This result should generally be output to the user as-is to indicate the content of the

API call. Additionally, if multiple API calls are involved, all relevant API call results must be merged and presented together to the user.

- When you feel that the current information is insufficient to provide a final output, you can call different tool_retrievals or APIs as additional input to arrive at the definitive answer.

- Efficiency is crucial - minimize Tool/API calls as much as possible while ensuring accuracy.

3. Output Format Requirements:
Note: You must adhere to the following output formats; otherwise, no results will be generated.

- When you determine that additional Tool/API calls are needed (Tool call format: API/tool_retrievals): $< think >$Thought process$< /think >< tool\_call >${"name": "tool_name", "arguments": {"param": "value"/["value"]}}$< /tool\_call >$

- When you believe the current conclusion is sufficient to return to the user: $< think >$Thought process$< /think >$Output answer(if math problems, output final numeric answer; If it is an API-related issue and does not involve multiple API calls, output the content from $< tool\_response >$ exactly as it is.)

---

**Query:** What is the result of 154+843?

**Recalled Memories:**

**Question 1**
Calculate the result of 222+333
**Answer**
555

<think>
The user's question is about calculating the sum of two numbers, 154 and 843. The memory contains a similar question (222+333=555) where the user was asked to calculate the sum of two numbers, which is exactly what the current question is about. Therefore, the answer can be directly used from the memory.
</think>

**Final Answer:** 555.                                                          ✗

**Memories**

**Question 1**                                                                  delete
Calculate the result of 222+333
**Answer**
555

| 👤 User | 🧠 Memory Module | 🤖 JoyAgents-R1 |

Figure 4: An example of dynamic updating of the agent memory module.

## A.7 Case analysis

**Dynamic update of the agent memory module.** Fig. 4 illustrates the dynamic updating mechanism of the math-agent's memory. When a user inputs an addition problem, the agent searches its memory repository and retrieves a similar problem whose numerical values differ but share the same structural pattern. However, due to interference from the recalled memory, the agent erroneously assumes the answer can be directly extracted from the memory without re-evaluating the new numerical values, resulting in an incorrect response. This erroneous memory entry is subsequently penalized through

Figure 5: The example of the QA agent in a question-answering task.

a scoring mechanism that reduces its retention probability in the memory buffer. Over successive iterations, such problematic memories are purged from the memory base, while memories that consistently contribute to accurate solutions are preserved. Eventually, the memory buffer converges to a stable state where it predominantly retains entries that enhance the agent's problem-solving reliability.

**More results for multiple tasks.** Fig. 5 ~Fig. 7 depict the reasoning processes of agents across distinct tasks, illustrating the hierarchical architecture and modular interactions within the system. Fig. 5 outlines the QA Agent's reasoning process: upon receiving a user query, the master agent routes the task to the QA Agent, which first invokes the retrieval tool to fetch relevant knowledge entries from the external knowledge base. The retrieved information is then processed by the summary tool to generate a concise answer, which is returned to the master agent for final delivery. Fig. 6 details the e-commerce function-call agent's operation. Internally, the agent begins by deploying the intention tool to parse user intent, followed by the retrieval tool to fetch intent-aligned APIs. Among the retrieved options, the agent selects the search_order_code API, executes it, validates the response, and relays the result back to the master agent. Fig. 7 presents a complete reasoning process from query to response. When a user inquires about the annual guarantee deposit for opening a furniture store on an e-commerce platform, the master agent first queries its memory buffer and identifies a similar historical case. This triggers the invocation of the e-commerce function-call agent. Within this sub-agent, the memory buffer is again consulted, recalling a prior solution that involved sequential use of the intention tool, retrieval tool, and check_shop_expenses API. The agent replicates this process. It analyzes user intent, retrieves contextually relevant APIs, selects the check_shop_expenses API based on semantic matching, and executes it to obtain the deposit amount (*e.g.*, "2000 yuan/year"). After validating the API response's relevance, the sub-agent returns the result to the master agent, which forwards it to the user. In summary, these cases underscore that task resolution necessitates coordination between the master agent and at least one specialized sub-agent, with memory buffers playing a pivotal role in both retrieving prior solutions and optimizing future responses through iterative learning.

**Query:** Hello, I have an order here. The buyer is Mark, and it was placed 1 hour and 10 minutes ago. I would like to confirm through which channel this order was placed.

**Recalled Memories: None**

E-commerce Function-call Agent

**<think>**
The user's query is about determining the total deposit required for a small furniture store run by an individual on [e-commerce platform A] for one year. ...... Therefore, the tool_intention tool should be called to further summarize and clarify the user's intention. This will help in identifying the exact deposit structure and the appropriate tool to retrieve the required information.
**</think>**

**<tool_call>**
{"name": "**tool_intention**", "arguments": {}}
**</tool_call>**

**<tool_response>**
Wants to confirm the channel through which the order placed by buyer Mark was made.
**</tool_response >**

**<think>**
The user's intention is to confirm the channel through which the order placed by buyer **Mark** was made. The reference information provided in <memory> does not contain any similar examples or dynamic programming strategies that could directly address this issue ...... To resolve this, it is necessary to call the tool_retrievals tool to further search for a suitable tool that can provide the required information about the order channel.
**</think>**

**<tool_call>**
{"name": "**tool_retrievals**", "arguments": {"intention": ["**Wants to confirm the channel through which the order placed by buyer Mark was made.**"]}}
**</tool_call>**

**<tool_response>**
{'name': 'search_order_code', 'description': 'Query the order details according to the information provided by the user, such as order number, order time, order account, customer name, product name, ......', 'parameters': {'type': 'object', 'properties': {'order_id': {'description': 'Order ID, Order, Unique identifier of order, multi-digit number.'}, 'consumer_name': {'description': 'Customer Name, User Name, Buyer Name, Recipient, Consignee, Customer, Client, etc.'}, ......}}
**</tool_response >**

**<think>**
The user's intention is to confirm the channel through which the order placed by buyer **Mark** was made ...... Therefore, the reasoning process concludes that calling the 'search_order_code' tool with the specified parameters is the appropriate solution to resolve the user's issue.
**</think>**

**<tool_call>**
{"name": "**search_order_code**", "arguments": {"consumer_name": ["**Mark**"], "time_context": ["1 hour and 10 minutes ago"], "search_keys": ["order source"]}}
**</tool_call>**

**<tool_response>**
Pass calling {'api_name': 'search_order_code', 'parameter': {'consumer_name': ['Mark'], 'search_keys': ['order source'], 'time_context': ['1 hour and 10 minutes ago']}}, the API call was successful.
Results are as follows: We have retrieved the relevant order information for you.
**</tool_response >**

**<think>**
The user's intention is to confirm the channel through which the order placed by buyer Mark was made. The historical dialogue indicates that the assistant has already identified this intention and called the tool 'search_order_code' to retrieve the relevant information, ....., as it aligns with the user's intention and the tool's successful execution.
**</think>**

**Pass calling {'api_name': 'search_order_code', 'parameter': {'consumer_name': ['Mark'], 'search_keys': ['order source'], 'time_context': ['1 hour and 10 minutes ago']}}, the API call was successful.**
**Results are as follows: We have retrieved the relevant order information for you.**

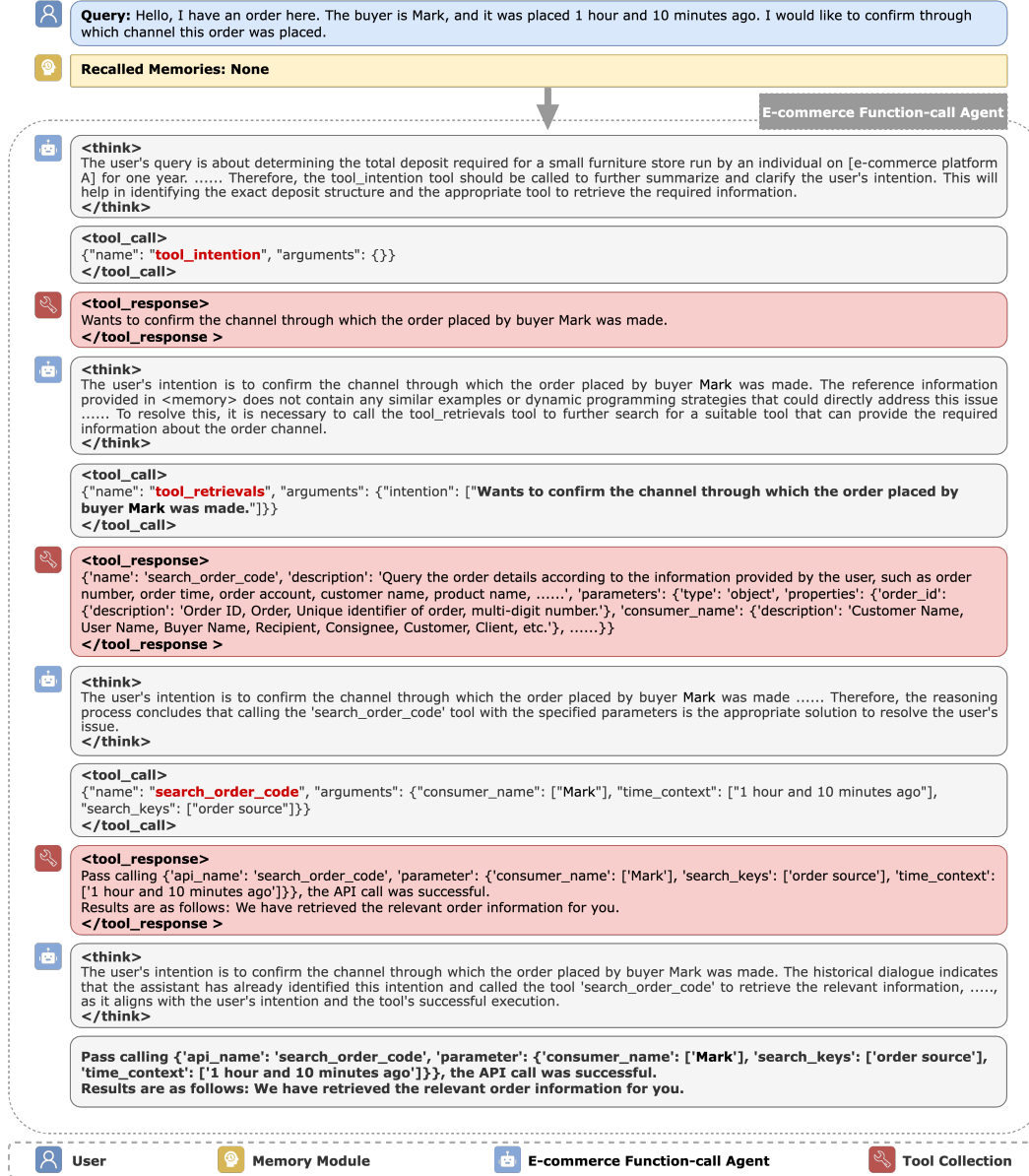| | | | |
|---|---|---|---|
| User | Memory Module | E-commerce Function-call Agent | Tool Collection |

Figure 6: The example of the function-call agent in an e-commerce task.
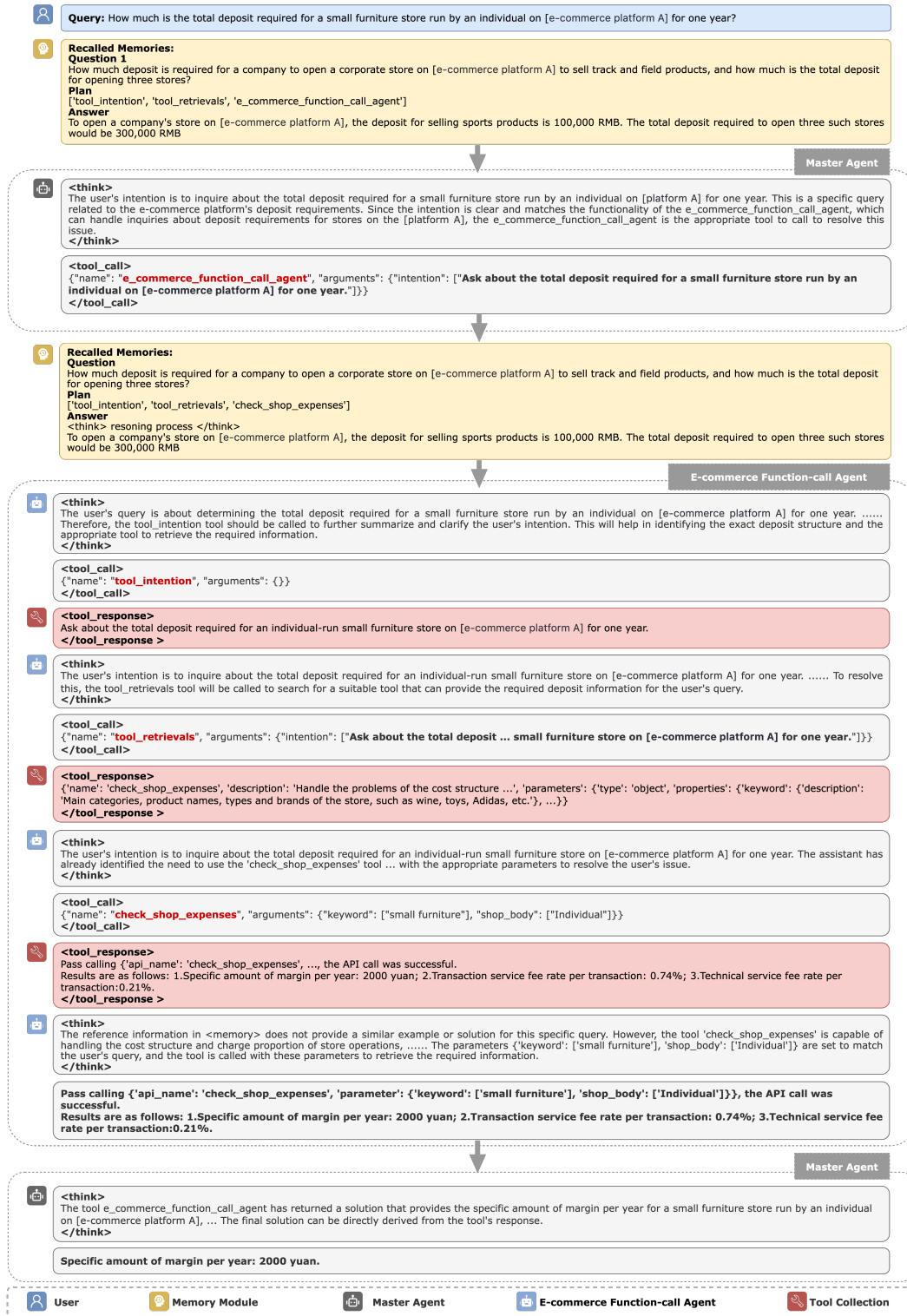
Figure 7: The example of a complete reasoning trajectory in an e-commerce task.