# Generalizing Experience for Language Agents with Hierarchical MetaFlows

Shengda Fan<sup>1</sup>, Xin Cong<sup>2\*</sup>, Zhong Zhang<sup>3</sup>, Yuepeng Fu<sup>3</sup>, Yesai Wu<sup>3</sup> Hao Wang, Xinyu Zhang, Enrui Hu, Yankai Lin<sup>1\*</sup>

<sup>1</sup>Gaoling School of Artificial Intelligence, Renmin University of China <sup>2</sup>Department of Statistics and Data Science, Tsinghua University <sup>3</sup>Department of Computer Science and Technology, Tsinghua University {fanshengda,yankailin}@ruc.edu.cn congxin1995@mail.tsinghua.edu.cn

#### Abstract

Recent efforts to employ large language models (LLMs) as agents have demonstrated promising results in a wide range of multi-step agent tasks. However, existing agents lack an effective experience reuse approach to leverage historical completed tasks. In this paper, we propose a novel experience reuse framework MetaFlowLLM, which constructs a hierarchical experience tree from historically completed tasks. Each node in this experience tree is presented as a MetaFlow which contains static execution workflow and subtask required by agents to complete dynamically. Then, we propose a Hierarchical MetaFlow Merging algorithm to construct the hierarchical experience tree. When accomplishing a new task, MetaFlowLLM can first retrieve the most relevant MetaFlow node from the experience tree and then execute it accordingly. To effectively generate valid MetaFlows from historical data, we further propose a reinforcement learning pipeline to train the MetaFlowGen. Extensive experimental results on AppWorld and WorkBench demonstrate that integrating with MetaFlowLLM, existing agents (e.g., ReAct, Reflexion) can gain substantial performance improvement with reducing execution costs. Notably, MetaFlowLLM achieves an average success rate improvement of 32.3% on AppWorld and 6.2% on WorkBench, respectively. The code is available at https://github.com/RUCBM/MetaFlowLLM.

# 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities in a wide range of multi-step agent tasks, including interactive coding [39, 32], web browsing [40, 46], and embodied housework [30, 20]. Existing work leveraging LLMs as autonomous agents that can iteratively take actions based on the execution feedback, enabling dynamic adaptation throughout the task execution process [41, 29]. Despite its effectiveness, when handling a large volume of tasks, earlier work tends to deal with each task independently, lacking the ability to summarize prior completed tasks to accumulate experience and enhance future task execution. As a result, the inability to learn from prior experiences results in the repetition of unnecessary steps (e.g., redundant trial-and-error) even when performing similar or identical tasks, thereby leading to ineffective and inefficient execution.

To address this limitation, two primary approaches are commonly employed: (1) **Trajectory-based Experience**, which involves constructing a database of historical task completion trajectories and retrieving the most relevant examples as in-context demonstrations to assist in executing new tasks [3, 45]; and (2) **Guideline-based Experience**, which involves summarizing prior experiences into

<sup>\*</sup> Corresponding author.

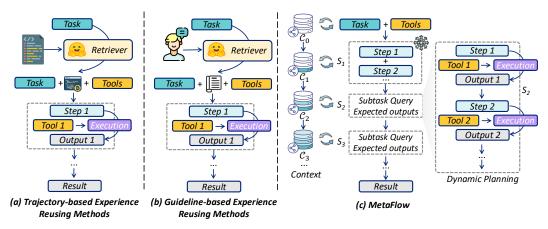


Figure 1: Illustration of three experience reuse approaches.

abstract natural language guidelines and using these guidelines to inform the execution of future tasks [44, 12]. However, trajectory-based experience often introduces fine-grained but task-irrelevant details, as even ostensibly similar tasks may exhibit divergent execution pathways, leading to reduced robustness when task contexts or environmental conditions change. Conversely, guideline-based experience typically presents experience in a coarse-grained form, which may omit critical procedural nuances, lacking the precision, structure, and robustness needed for reliable execution.

Humans do not learn experience by directly reusing specific procedures or merely summarizing abstract guidelines. A more effective approach is to construct an experience hierarchy based on the similarity between tasks: the more similar the tasks are, the more concrete procedures can be directly reused; the less similar they are, the more dynamic execution is needed [4, 7]. For instance, consider the following two tasks: (1) Send an email to attendees of the November 30 event with the title 'Remember to attend this event.' (2) Remind attendees of the December 1 event via email titled 'Remember to attend this event.' As humans, we can recognize that many steps in these two tasks can be reused, including searching for the event using calendar.search\_events, extracting participant emails, and sending emails via the email.send\_email tool. For those differences (e.g., extracting the event date from the user query), adjustments can be made based on the specifics of each task. Motivated by this human ability, we propose MetaFlowLLM, a novel framework that builds a hierarchical experience tree from historical task execution data to improve the effectiveness and efficiency of LLMs. Specifically, as illustrated in panel (c) of Figure 1, each node in the experience tree represents an abstracted workflow named **MetaFlow**, which contains both concrete execution workflow (i.e., the static step) and subtasks (i.e., the dynamic step) that need LLMs to execute dynamically. The closer to the leaf nodes, the more static steps are in the MetaFlow; the closer to the root node, the more dynamic steps it contains. When confronted with a new task, MetaFlowLLM retrieves from the hierarchical experience tree to identify the most appopriate MetaFlow node and initiates its execution. For components corresponding to concrete workflows, it executes them directly. For components requiring dynamic execution, an autonomous agent is responsible for interpreting the subtask description and completing it accordingly.

Implementing MetaFlowLLM presents two key challenges: (1) how to abstract MetaFlows given arbitrary two task execution trajectories: considering that manually summarizing MetaFlows is both time-consuming and labor-intensive, and that directly prompting powerful LLMs such as GPT-40 or DeepSeek R1 does not yield satisfactory results, we propose a reinforcement learning pipeline based on verifiable rewards to train an effective MetaFlow generator, called MetaFlowGen. Specifically, we introduce two distinct rewards: *i*) a **correctness reward**, which quantifies the proportion of queries for which the generated MetaFlow yields correct results, and *ii*) an **efficiency reward**, which incentivizes the preservation of static steps while minimizing unnecessary invocations of LLMs. Leveraging these rewards, we initially construct a high-quality supervised fine-tuning dataset via rejection sampling to enable effective cold-start. Subsequently, we directly optimize the MetaFlow generator using reinforcement learning to maximize the defined rewards. (2) how to construct an experience tree from existing historical data: we draw inspiration from hierarchical clustering techniques [24, 25] and propose a Hierarchical MetaFlow Merging algorithm. This algorithm iteratively merges the

similar nodes using the MetaFlowGen until the tree is fully constructed. To ensure the usability of the experience tree, we incorporate a pruning algorithm based on the reward of the merged nodes, which helps maintain its efficiency and effectiveness.

To validate the effectiveness of our proposed method, we conduct extensive experiments on two distinct scenarios including AppWorld [32] and WorkBench [31]. Experimental results demonstrate that integrating MetaFlow into existing LLM-based agents leads to substantial performance improvements while simultaneously reducing execution costs. Notably, MetaFlowLLM achieves an average improvement of 32.3% on AppWorld and 6.2% on WorkBench, surpassing both trajectory-based and guideline-based experience reuse approaches.

Our contributions are as follows:

- We introduce MetaFlowLLM as a novel mechanism for experience reuse in LLM-based agents, enabling them to generalize structured workflows across similar tasks.
- We propose a reinforcement learning-based training pipeline for the MetaFlow generator to improve the MetaFlow generation capabilities of LLMs.
- Through comprehensive experiments on AppWorld [32] and WorkBench [31], we demonstrate that our method produces high-quality workflows that significantly improve LLM performance and reduce computational overhead.

# 2 Related Work

**LLM-based Agents** With advances in the instruction-following and reasoning capabilities of LLMs [1, 17, 38, 13], their applications have expanded beyond traditional natural language processing tasks such as question answering [27] and information extraction [10, 9]. Researchers have begun to use LLMs as autonomous agents for tasks such as code generation [26, 32], travel planning [37], web browsing [40, 46], and embodied housework [30, 20]. Currently, the two most representative LLM agent frameworks are ReAct [41] and Reflexion [29]. However, these two vanilla agents are incapable of utilizing historical data to form experience to promote future tasks. To address this limitation, we propose the MetaFlowLLM framework, which enables agents to accumulate, retrieve, and leverage past experiences to enhance reasoning, decision-making, and generalization across tasks.

Experience Reuse in LLM-based Agents Current approaches to experience reuse in LLM agents primarily follow two paradigms: (1) trajectory-based experience: direct utilization of complete or abstracted trajectories as contextual examples [3, 45, 34], which often suffer from overfitting to specific instances, and (2) guideline-based experience: manual or LLM-generated summarization of experiences into natural language instructions that are incorporated into the agent's prompt [44, 12]. For the former one, they retain excessive task-specific details and not sufficiently generalizable. For the latter one, the abstract guideline in natural language may lacking the precision, structure, and robustness needed for reliable execution. To address these limitations, we propose a novel experience reuse framework, MetaFlowLLM, that constructs a hierarchical experience tree covering experiences at different levels. Compared to AWM [34], which extracts reusable workflows from historical trajectories and still relies on prompt-based experience reuse driven by the LLM's reasoning capability, MetaFlowLLM explicitly models reusable workflows as structured graphs with static and dynamic nodes, and further optimizes their granularity through reinforcement learning.

Workflow Generation Process Automation [5] has long been a key goal in technology, aiming to automate repetitive tasks and improve operational efficiency. Traditional Robotic Process Automation relies on manually annotated workflows to abstract such tasks [19, 15, 35, 2, 11]. With the advent of LLMs, many efforts have focused on leveraging their capabilities to replace manual labor in workflow generation [42, 43, 16, 36, 21]. However, generating large-scale workflows with complex logical structures based solely on queries and available tool lists remains a significant challenge for current LLMs [8, 42], primarily due to the lack of intermediate execution feedback. To address this limitation, our work takes a more pragmatic approach by refraining from generating all specific execution steps. Instead, we abstract similar workflows into a unified MetaFlow, where shared steps are fixed and sample-specific components are delegated to sub-tasks handled by dynamic agents. In this way, MetaFlowLLM provides a balanced framework for automation—capturing common workflow patterns while maintaining the flexibility to adapt to diverse tasks.

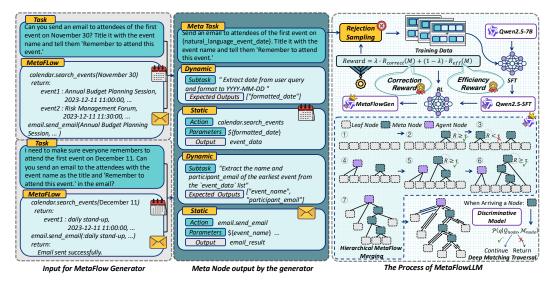


Figure 2: Illustration of the MetaFlowLLM framework. On the left, the task and corresponding MetaFlow (since the example in this figure is from the leaf node, it is referred to as a trajectory) serve as inputs to the generator. The middle section depicts the meta node (comprising the meta task and MetaFlow) generated by the MetaFlow generator. On the right, the process of MetaFlowLLM framework is outlined, including: 1) the training process, where a rejection sampling approach constructs the cold-start dataset, followed by SFT and RL to train the MetaFlowGen; and 2) the deployment process, which involves constructing the hierarchical experience tree and performing deep matching traversal to find the most appropriate MetaFlow node.

# 3 Methodology

In this section, we introduce the MetaFlowLLM framework, which constructs a hierarchical experience tree containing different level MetaFlows to leverage historical completed tasks. We first formally define the MetaFlow structure and its operational rules in Section 3.1. Then, we propose a reinforcement learning pipeline to derive a dedicated MetaFlowGen to generate MetaFlows, as detailed in Section 3.2. Finally, we present the hierarchical experience tree construction and deep matching traversal algorithms in Section 3.3.

# 3.1 MetaFlow Definition

As illustrated in part (c) of Figure 1, MetaFlow is a hybrid execution framework composed of static procedural steps and dynamic LLM-based agent steps, all operating over a shared and mutable context  $\mathcal{C}$ . Formally, a MetaFlow  $\mathcal{M} = \{s_1, s_2, \ldots, s_n\}$  is defined as a sequence of steps, each of which can be either a static step  $s_i^S$  or a dynamic step  $s_i^D$ . The shared context  $\mathcal{C}$  is initialized with the user task. As the MetaFlow progresses,  $\mathcal{C}$  is iteratively updated by both static steps and the outputs of online planning agents. In practice,  $\mathcal{C}$  may be an implicit variable scope (e.g., in Python), or an explicit key-value store tracking agent decisions and tool outputs.

At runtime, given a task q, MetaFlowLLM executes each step sequentially in the order defined by  $\mathcal{M}$ . For each step  $s_i$ , the framework checks the step type:

- If  $s_i$  is a static step, consisting of operations such as tool calls or code execution, it is directly executed, and the context is subsequently updated as  $C_i = s_i^S(C_{i-1})$ .
- If  $s_i$  is a dynamic agent step, MetaFlowLLM adapts the subtask of  $s_i$  based on the task q and the current context  $C_{i-1}$ , producing a state-specific prompt  $p_i$  as subtask prompt.  $p_i$  is then used to initialize an LLM-based agent A, such as ReAct [41] or Reflexion [29]. The LLM agent A interacts with the environment over multiple turns, each comprising an action and the subsequent observation. Formally, the agent A produces an action sequence  $\{a_i^1, a_i^2, \ldots, a_i^k\} = A(p_i)$ , where each action  $a_i^j$  is executed and the context is updated accordingly:

$$C_i^j = s_i^D(C_i^{j-1}, a_i^j), \quad C_i^0 = C_{i-1}.$$

The process continues until the agent determines that the subtask is complete. The final context after executing all actions is denoted as  $C_i = C_i^k$ .

# 3.2 Training MetaFlowGen with Reinforcement Learning

To generate effective MetaFlows, we adopt a reinforcement learning pipeline to train a MetaFlow generator named MetaFlowGen. As illustrated by the right part of Figure 2, the pipeline follows a two-stage approach: supervised fine-tuning (SFT) to initialize the model for cold-start, and rule-based reinforcement learning using GRPO [28] to enhance the MetaFlow generation capability.

As illustrated in the left part of Figure 2, the MetaFlow generator is designed to take two samples as input: each consisting of tasks with the corresponding MetaFlows and the relevant tool documentation. The goal is to generate a high-level MetaFlow  $\mathcal M$  that captures their common patterns through static steps and handles their different procedures using dynamic agents. Formally, the MetaFlow generation process can be expressed as:

$$\mathcal{M}, \mathcal{Q} = \pi_{\theta} \left( q_1, q_2, m_1, m_2, \mathcal{A} \right), \tag{1}$$

where  $\theta$  is the parameter of the MetaFlowGen,  $q_1,q_2$  are the input tasks respectively,  $m_1,m_2$  are their corresponding MetaFlows, and the tool documentation  $\mathcal A$  provides necessary background knowledge. As illustrated in the middle part of Figure 2. the output of the MetaFlowGen consists of the MetaFlow  $\mathcal M$  and its associated meta task  $\mathcal Q$ .

Supervised Fine-Tuning for Cold-Start. Following previous work [14], we use SFT as a cold-start method. The SFT training dataset, denoted as  $D_{\rm SFT} = \{(q_1,q_2,m_1,m_2,\mathcal{M},\mathcal{Q})^{(i)}\}_{i=1}^{|D_{\rm SFT}|}$ , is constructed through rejection sampling [6]. Specifically, we maintain a list of golden nodes, initialized with the samples from the training set (referred to as leaf nodes). In each iteration, we sample two golden nodes and pair them with manually crafted in-context learning samples to form an input prompt. This prompt is then used to generate the corresponding meta task and MetaFlow via an LLM as a MetaFlow generator. We execute this MetaFlow on all its leaf nodes, and if it passes all the test cases associated with these nodes (i.e., its correctness reward equals 1.0 defined in Equation 3), it is deemed correct and added to the golden nodes list. This iterative process continues until the dataset reaches the desired scale. It effectively expands both the dataset size and the complexity of the hierarchical structures, facilitating the subsequent construction of the hierarchical experience tree.

Formally, the objective function for supervised fine-tuning is defined as:

$$\mathcal{L}_{SFT}(\theta) = -\mathbb{E}_{(q_1, q_2, m_1, m_2, \mathcal{M}, \mathcal{Q}) \sim D_{SFT}} \left[ \log \pi_{\theta}(\mathcal{M}, \mathcal{Q} \mid q_1, q_2, m_1, m_2, \mathcal{A}) \right]. \tag{2}$$

**Rule-based Reinforcement Learning.** After completing the SFT phase, we further enhance the model using the reinforcement learning approach. In this phase, the training dataset is constructed by sampling from the list of golden nodes produced during SFT. To guide the learning process, we design a reward function that integrates two key objectives: correctness and efficiency.

• **Correctness Reward.** This term evaluates whether initializing the MetaFlow correctly solves the leaf nodes' tasks. The correctness reward is defined as the ratio of correctly solved leaf nodes:

$$R_{\text{correct}}(\mathcal{M}) = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(\text{execute}(q_i, \mathcal{M}) = \text{expected}_i), \tag{3}$$

where N is the number of leaf nodes, and  $\mathbb{I}$  is the indicator function that checks whether the generated result matches the expected output.

• Efficiency Reward. This component evaluates the preservation of static actions and runtime efficiency. The efficiency reward is defined as the ratio of static operations to the minimum number of static operations between the two input workflows:

$$R_{\rm eff}(\mathcal{M}) = \frac{{\tt NumStaticOps}(\mathcal{M})}{\min({\tt NumStaticOps}(m_1), {\tt NumStaticOps}(m_2))}, \tag{4}$$

This encourages workflows that maximize static step reuse and minimize redundant dynamic calls, improving both reusability and efficiency.

# Algorithm 1 Hierarchical MetaFlow Merging

```
Require: Leaf nodes \mathcal{D} = \{(q_i, m_i)\}, distance metric d(\cdot, \cdot), MetaFlowGen \theta, correctness threshold \tau
Ensure: Hierarchical experience Tree \mathcal{T}
 1: Initialize: Initialize the set of nodes \mathcal{N} \leftarrow \mathcal{D}
2: while |\mathcal{N}| > 1 do
3:
       Select the most similar node pair: From \mathcal{N}, select two most similar nodes n_1, n_2 based on the distance
        metric d(\cdot, \cdot)
        Generate candidate MetaFlow: Merge n_1 and n_2 using the MetaFlowGen \theta to produce a candidate
        MetaFlow \mathcal{M}
        Evaluate its correctness: Compute the correctness reward R_{\text{correct}} for \mathcal{M}
6:
        if R_{\rm correct} > = \tau then
7:
           Use \mathcal{M} as the new parent node to merge n_1 and n_2
8:
           Remove n_1, n_2 from \mathcal{N} and add \mathcal{M} to \mathcal{N}
9:
10:
           Use pure-agent node as the new parent node to merge n_1 and n_2
11:
           Remove n_1, n_2 from \mathcal{N}
12:
        end if
13: end while
14: Final merge: If there are more than one pure-agent nodes, merge them into a single top-level pure-agent
    node
15: return The complete experience Tree \mathcal{T}
```

The total reward is the weighted sum of the above two rewards<sup>2</sup>, guiding the generator to produce MetaFlows that are both accurate and efficient.

$$R(\mathcal{M}, \mathcal{Q}) = \lambda \cdot R_{\text{correct}}(\mathcal{M}) + (1 - \lambda) \cdot R_{\text{eff}}(\mathcal{M}), \tag{5}$$

where  $\lambda$  is a scaling factor for controlling the trade-off between accuracy and efficiency.

# 3.3 Hierarchical Experience Tree Construction and Utilization

For efficient inference, we begin by constructing a hierarchical experience tree from the leaf nodes in the historical completed tasks. This approach is detailed in Algorithm 1. As shown in the right part of Figure 2, the set of nodes to be merged, denoted as  $\mathcal{N}$ , is initialized as the set of leaf nodes  $\mathcal{D}$ .  $\mathcal{N}$  is progressively merged and updated based on semantic similarity. Each merger results in an intermediate meta node  $\mathcal{M}$ . To ensure the effectiveness of the experience tree, we calculate the correctness reward of  $\mathcal{M}$  using Equation 3. If the reward exceeds a threshold  $\tau$ , the meta node is added to  $\mathcal{N}$ . Otherwise, the two nodes are directly merged into a pure agent node.

At inference time, the goal is to efficiently identify the most appropriate MetaFlow to solve the task. As shown in Algorithm 2 and Figure 2, the deep matching traversal process begins by matching the test task  $q_{\text{test}}$  with the experience tree  $\mathcal{T}$ , starting from the root node and moving downwards. The node-level matching is performed using a discriminative model  $\mathcal{P}(q|\mathcal{Q},\mathcal{M})$ , which prompts an LLM to determine whether the task q can be classified as an instance of a  $\mathcal{Q}$  and can be solved by the  $\mathcal{M}$ . The search continues recursively through child nodes until the deepest matching node is found.

# 4 Experiments

#### 4.1 Experimental Settings

**Datasets and Metrics.** We conduct experiments on two representative agent datasets: AppWorld [31] and WorkBench [32]. AppWorld serves as a representative dataset for interactive coding tasks, while WorkBench serves as a representative dataset for tool learning tasks in workplace settings. The evaluation metric in AppWorld is the **Task Goal Completion (TGC) Rate**, which is defined as the percentage of tasks where the agent passes all evaluation tests. We use two evaluation metrics in WorkBench: (1) **Accuracy**, which is defined as the percentage of tasks where the outcome from the agent's actions matches the expected outcome of the ground truth; and (2) **Side Effects**, which quantitatively evaluate unintended changes caused by tool calls. We conduct experiments under two

<sup>&</sup>lt;sup>2</sup>We also incorporate a format-based penalty [14] in the reward design: a reward of -1 is assigned when the  $\mathcal{M}$  fails to meet the required format constraints.

Table 1: TGC rate comparison of various models with different methods on AppWorld(	rate comparison of various models with different methods on AppWorld(%	)).
--	--	-----

Model	Agent Type	Configuration				
	<b>g, F</b> -	Base	w/ Traj.	w/ Guideline	w/ MetaFLow	
Qwen2.5-7B	Reflexion	7.6	32.7	13.5	37.4	
	ReAct	6.4	37.4	17.0	43.9	
Qwen2.5-32B	Reflexion	12.9	33.9	22.8	43.9	
	ReAct	22.2	49.1	38.6	50.9	
GPT-4o-mini	Reflexion	13.5	37.4	20.5	45.6	
	ReAct	7.0	22.2	7.0	35.7	

conditions: one where only the required tools are provided to LLMs, and another where all the tools are available. The dataset statistics can be found in Appendix D.

**Baselines.** To evaluate the performance of MetaFlow, we consider four base LLMs during the inference stage: two state-of-the-art open-source models, Qwen2.5-7B-Instruct [38] and Qwen2.5-32B-Instruct, and two proprietary models, GPT-4o-mini [18] and GPT-4o. In order to assess the generalizability of MetaFlow, we consider two agent frameworks as the backbone for dynamic nodes: (1) ReAct [41], which combines reasoning and acting, and (2) Reflexion [29], which enables the agent to leverage environmental feedback to reflect on mistakes and perform subsequent actions. We consider the following three baselines: (1) LLM agent without experience reuse, (2) Trajectory-based Experience reuse method (Denoted as *Traj.*), and (3) Guideline-based Experience reuse method (Denoted as *Guideline*). Further implementation details of the trajectory-based experience reuse and guideline-based experience methods can be found in Appendix E.

**Implementation Details.** We set  $\tau$  to 1.0 in all experiments to ensure the quality of the experience tree. The value of  $\lambda$  is set to 0.7. We use the all-MiniLM-L6-v2<sup>3</sup> model to encode the task, and the cosine similarity of the task embeddings is utilized as the distance metric.

In the SFT stage, we fine-tune the MetaFlowGen on Qwen2.5-7B-Instruct for 3 epochs using the AdamW optimizer [23] and a linear learning rate scheduler with a peak learning rate of  $2\times 10^{-5}$ . Each mini-batch contains 32 examples, and the maximum sequence length is set as 8,192 tokens. In RL stage, we adopt TRL [33] as our training framework. We set the training epochs to 2, batch size to 28, learning rate to  $1\times 10^{-6}$ , KL coefficient to 0 [22], rollout number to 14. All experiments are conducted on 8 NVIDIA A800 40G GPUs.

#### 4.2 Main Results

**Results on AppWorld.** Following [32], we evaluate the performance of LLMs using the ReAct and Reflexion prompting methods. The results are placed in Table 1, from which we derive that:

- 1. The improvement of the guideline methods on the challenging AppWorld is significantly smaller than that of trajectory-based methods. This suggests that, when faced with difficult tasks, the natural language guideline method may omit critical procedural details, lacking the precision, structure, and robustness required for reliable execution.
- 2. The proposed MetaFlowLLM method can effectively adapt to both the Reflexion and ReAct agent frameworks, leading to substantial performance improvements. Notably, MetaFlow increases TGC by 33.1% and 31.6% compared to the base Reflexion and ReAct agents, consistently outperforming the trajectory-based and guideline-based experience reuse methods. The results demonstrate the effectiveness of breaking down complex tasks into distinct sub-tasks, with clear and reusable transition logic between them. This decomposition allows to reuse static code and agent's plan from historically similar tasks, thereby increasing both execution stability and reliability.

**Results on WorkBench.** Following [31], we assess the performance of LLM agents using the ReAct [41] framework. Since the WorkBench dataset lacks pre-defined "golden" trajectories, we

<sup>3</sup>https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2

Table 2: Accuracy and side effects comparison of various models with different methods on Work-Bench(%). The best results are marked in **bold** and the second-best results are marked with <u>underline</u>. Abbreviations are defined as follows. CRM: Customer Relationship Manager, PM: Project Management, MD: Multi-Domain.

Model	Anal	ytics	Cale	ndar	CR	M	Em	ail	PN	M	M	D	Av	g
1,10401	%acc↑	%se↓	%acc↑	%se↓	%acc↑	%se↓	%acc↑	%se↓	%acc↑	%se↓	%acc↑	%se↓	%acc↑	%se↓
GPT-4o	30.8	53.8	53.0	19.7	30.9	14.5	32.8	25.9	8.16	0.0	16.3	48.9	27.2	33.1
w/ Traj.	48.5	47.4	48.5	10.6	54.5	10.9	41.3	25.9	20.4	4.13	25.1	55.1	35.3	32.7
w/ Guideline	34.6	44.9	53.0	21.2	40.0	23.6	31.0	27.6	8.16	0.0	15.6	52.4	28.5	34.2
w/ MetaFlow	44.9	37.2	63.7	12.1	60.0	14.5	46.6	24.1	24.5	6.12	24.5	40.8	40.8	26.9
Qwen2.5-7B	18.2	60.3	27.6	35.5	16.0	30.5	9.31	43.8	8.57	6.12	9.39	53.9	14.3	43.0
w/ Traj.	30.8	46.2	18.2	48.5	18.2	32.7	15.5	46.6	14.3	12.2	8.16	57.1	16.3	44.8
w/ Guideline	21.8	51.3	37.9	47.0	20.0	52.7	15.5	55.2	10.2	6.12	8.16	60.5	17.4	49.5
w/ MetaFlow	<u>26.9</u>	43.6	25.8	50.0	23.6	34.5	10.3	55.2	<u>14.3</u>	18.4	9.52	57.8	<u>17.2</u>	46.8
Qwen2.5-32B	6.41	88.5	57.6	24.2	12.7	43.6	27.6	24.1	6.12	0.0	12.2	61.9	19.2	47.2
w/ Traj.	25.6	70.5	37.9	40.9	41.8	20.0	8.62	58.6	20.4	12.2	15.0	60.5	23.2	49.0
w/ Guideline	19.2	79.5	59.1	21.2	27.3	23.6	31.0	36.2	10.2	2.0	15.0	57.8	25.2	43.3
w/ MetaFlow	<u>24.4</u>	69.2	59.1	19.7	27.3	23.6	37.9	<u>32.8</u>	12.2	4.1	18.4	55.8	28.3	40.4
GPT-4o-mini	11.5	44.9	48.5	30.3	21.8	36.4	24.1	58.6	8.16	6.12	12.2	56.5	19.6	43.0
w/ Traj.	12.8	44.9	39.4	25.8	27.3	36.4	20.7	48.3	12.2	12.2	14.3	57.1	19.9	41.9
w/ Guideline	12.8	43.6	47.0	37.9	27.3	30.9	31.0	51.7	12.2	4.10	12.9	60.5	21.9	43.5
w/ MetaFlow	19.2	30.8	53.0	15.2	16.4	23.6	31.0	46.6	16.3	4.08	22.4	51.7	26.1	33.6

construct an *experience pool* via rejection sampling from previously collected trajectories. The results are shown in Table 2, from which we draw three observations:

- 1. All three methods result in improvements in the model's overall performance. However, the MetaFlow method demonstrates the most significant enhancement. Specifically, when employing MetaFlow, the average increase in accuracy is 6.16%, accompanied by a 4.13% reduction in side effects. Furthermore, when using GPT-40-mini as the backbone LLM, MetaFlow outperforms both trajectory-based and guideline-based experience methods across the six domains and evaluation metrics, with only one exception. This demonstrates the effectiveness of MetaFlowLLM in tool-calling tasks, as it can reuse experience more efficiently.
- 2. MetaFlow shows greater performance improvements with more powerful LLMs. For example, the accuracy of Qwen2.5-32B-Instruct improves by 9.1%, while Qwen2.5-7B-Instruct only increases by 2.9%, which is slightly lower than the performance of the guideline method at the same scale. Manual inspection reveals two key reasons: (1) low-capacity models exhibit insufficient discriminative capability to identify appropriate MetaFlows; (2) their weak instruction-following performance impedes the successful execution of MetaFlow sub-tasks.

#### 4.3 Analysis

Scalability Analysis. Theoretically, both MetaFlowGen training and tree construction are performed offline, while online retrieval operates with sublinear time complexity relative to the number of nodes. The most computationally intensive operation—correctness reward computation—formally scales as  $\mathcal{O}(N)$  with the number of leaf nodes but can be fully parallelized across GPUs. Empirically, the runtime of both RL training and tree construction grows approximately linearly with the number of experiences (see Figure 3).

MetaFlow can reduce the impact of noisy context from more available tools. Selecting inappropriate MetaFlow nodes can hinder performance. To investigate the significance of the discriminative model  $\mathcal{P}(q|\mathcal{Q},\mathcal{M})$ , we evaluate MetaFlowLLM's performance under hierarchical level shifts. For each test task, we employ GPT-40 as the discriminative model to identify appropriate MetaFlow nodes. Subsequently, we evaluate performance deviations when accessing either parent nodes (with negative level shifts) or child nodes (with positive level shifts). As shown in Figure 4, both upward and downward shifts lead to degradation in model accuracy,

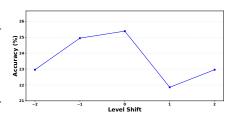
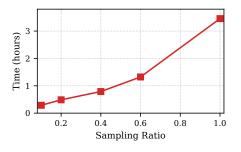
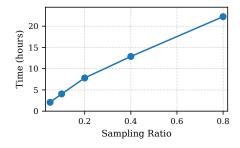


Figure 4: Comparison of accuracy with varying level shifts on WorkBench.





(a) Tree-building on WorkBench.

(b) RL training on AppWorld.

Figure 3: **Empirical runtime scaling of MetaFlowLLM.** Both curves exhibit approximately linear scaling as the sampling ratio increases.

indicating that the originally selected node is better than others. This demonstrates the importance of the discriminative model in the inference.

In our main experiments, we only provide the necessary tools to the ReAct agent. In this study, we evaluate the performance of MetaFlowLLM in a more real-world setting where all 26 tools are provided. As shown in Table 3, MetaFlowLLM results in a more considerable improvement in accuracy, while concurrently reducing side effects. For example, Qwen2.5-7B-Instruct achieves an accuracy improvement from 9.05% to 15.7%. This indicates that MetaFlowLLM mitigates the effect of noisy context by predefining tool usage in static steps and clearer planning

Table 3: Accuracy and side effects comparison with all tools on WorkBench (%).

Model	Accuracy (%)	Side Effect (%)
Qwen2.5-7B	9.05	47.9
w/ MetaFlow	15.7	41.7
Qwen2.5-32B	13.8	52.3
w/ MetaFlow	23.2	47.7

predefining tool usage in static steps and clearer planning in dynamic steps, thereby reducing the complexity of tool selection.

The MetaFlowGen demonstrates superior performance over proprietary models. To validate the effectiveness of the proposed training pipeline, we conduct comparative experiments involving GPT-40, DeepSeek-R1, and our MetaFlowGen. We provided three in-context learning examples for GPT-40 and DeepSeek-R1, while the trained models utilize a zero-shot prompt. The experimental results (Table 4) reveal two findings: (1) The trained MetaFlowGen achieves significantly higher performance than both GPT-40 and DeepSeek-R1 across both

Table 4: Comparison between various MetaFlow generators on AppWorld.

Model	Correct	Efficiency
GPT-40	0.160	0.60
DeepSeek-R1	0.054	0.41
MetaFlowGen	0.360	0.99
w/o RL	0.160	0.80
w/o SFT	-0.380	-0.14

metrics; (2) Ablation studies demonstrate consistent performance degradation when either SFT or RL is removed, confirming the essential contribution of both proposed training stages. We also provide the reward curves during RL training and more detailed analysis in Appendix B.

The proposed RL framework exhibits superior sample efficiency compared to the SFT baseline. To evaluate the cost-effectiveness of our reinforcement learning pipeline, we vary the sampling ratio of the RL training set and analyzed the resulting correctness and efficiency metrics, as shown in Figure 5. We observe a consistent improvement in both metrics as the sampling ratio increases. Notably, the model attains near-optimal performance at a full sampling ratio (1.0), achieving a peak efficiency of 0.991. Although the correctness slightly declines at this point, the overall gain in efficiency highlights the inherent trade-off between accuracy and computational cost—an expected characteristic in RL optimization. These results demonstrate that the RL training method can achieve competitive performance with substantially fewer training samples, underscoring its strong sample efficiency.

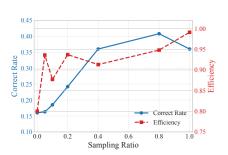


Figure 5: Sample efficiency of the proposed RL framework.

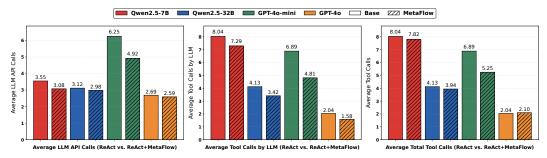


Figure 6: Efficiency comparison between ReAct and ReAct+MetaFlow, based on the average number of LLM API calls and tool calls.

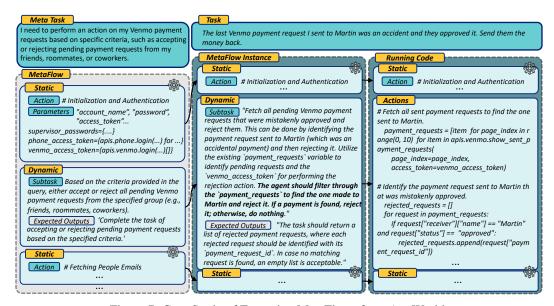


Figure 7: Case Study of Executing MetaFlows from AppWorld.

MetaFlowLLM can enhance the efficiency of sequential decision-making. We evaluate the performance of different models using ReAct and ReAct+MetaFlow based on three key metrics: the number of LLM API calls, the average number of tool calls by LLMs, and the average number of tool calls overall. As shown in Figure 6, all models demonstrate a consistent reduction in both the number of LLM API calls and the number of tool calls by LLMs when MetaFlow is applied. This indicates that the introduced static steps effectively address subproblems. In addition, we find that MetaFlowLLM helps reduce the overall number of tool calls, particularly when the base LLMs are relatively weak. Manual inspection indicates that the elevated tool call frequency is primarily due to the model's dependence on repeated trial and error during execution. More efficiency comparison results can be found in Appendix H.

Case Study of Executing MetaFlows. Figure 7 presents a MetaFlow execution example in AppWorld. The meta task and its corresponding MetaFlow are shown in the left part of the figure. Static steps are executed as-is, while dynamic steps are adapted according to the input task and context (transition from the left to the center part of Figure 7). The resulting subtasks are then used to prompt the agent, which generates and executes code accordingly (center to right part of Figure 7).

# 5 Conclusion

In this paper, we present MetaFlowLLM, a novel framework that enables structured experience reuse for LLM-based agents through a hierarchical experience tree. By abstracting past task trajectories into a tree of MetaFlows, our approach facilitates effective and efficient task execution for new-coming tasks. In future work, we plan to extend MetaFlowLLM to multimodal domains and beyond.

# Acknowledgements

This work was supported by the National Key R&D Program of China(No.2024YFC3306500), Beijing Nova Program (No. 20240484568), the Postdoctoral Fellowship Program of CPSF (Grant No. GZB20230343 and Grant No. GZC20240831) and the China Postdoctoral Science Foundation (Grant No. 2023M741945 and Grant No. 2025M771586).

# 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] Simone Agostinelli, Andrea Marrella, and Massimo Mecella. Towards intelligent robotic process automation for bpmers. *arXiv preprint arXiv:2001.00804*, 2020.
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
- [4] Michelene TH Chi, Paul J Feltovich, and Robert Glaser. Categorization and representation of physics problems by experts and novices. *Cognitive science*, 5(2):121–152, 1981.
- [5] Andrzej Cichocki, Helal A Ansari, Marek Rusinkiewicz, and Darrell Woelk. *Workflow and process automation: concepts and technology*, volume 432. Springer Science & Business Media, 1997.
- [6] Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, SHUM KaShun, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *Transactions on Machine Learning Research*.
- [7] Maria K Eckstein and Anne GE Collins. How the mind creates structure: Hierarchical learning of action sequences. In *CogSci... Annual Conference of the Cognitive Science Society. Cognitive Science Society (US). Conference*, volume 43, page 618, 2021.
- [8] Shengda Fan, Xin Cong, Yuepeng Fu, Zhong Zhang, Shuyan Zhang, Yuanwei Liu, Yesai Wu, Yankai Lin, Zhiyuan Liu, and Maosong Sun. WorkflowLLM: Enhancing workflow orchestration capability of large language models. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [9] Shengda Fan, Shasha Mo, and Jianwei Niu. Boosting document-level relation extraction by mining and injecting logical rules. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10311–10323, 2022.
- [10] Shengda Fan, Yanting Wang, Shasha Mo, and Jianwei Niu. Logicst: A logical self-training framework for document-level relation extraction with incomplete annotations. In *Proceedings* of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 5496– 5510, 2024.
- [11] Deborah Ferreira, Julia Rozanova, Krishna Dubba, Dell Zhang, and Andre Freitas. On the evaluation of intelligent process automation. *arXiv* preprint arXiv:2001.02639, 2020.
- [12] Yao Fu, Dong-Ki Kim, Jaekyeom Kim, Sungryull Sohn, Lajanugen Logeswaran, Kyunghoon Bae, and Honglak Lee. Autoguide: Automated generation and selection of context-aware guidelines for large language model agents. *arXiv preprint arXiv:2403.08978*, 2024.
- [13] 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.
- [14] 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.

- [15] Peter Hofmann, Caroline Samp, and Nils Urbach. Robotic process automation. *Electronic markets*, 30(1):99–106, 2020.
- [16] Tian Huang, Chun Yu, Weinan Shi, Zijian Peng, David Yang, Weiqi Sun, and Yuanchun Shi. Promptrpa: Generating robotic process automation on smartphones from textual prompts. *arXiv* preprint arXiv:2404.02475, 2024.
- [17] Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.
- [18] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv* preprint arXiv:2410.21276, 2024.
- [19] Lucija Ivančić, Dalia Suša Vugec, and Vesna Bosilj Vukšić. Robotic process automation: systematic literature review. In *Business Process Management: Blockchain and Central and Eastern Europe Forum: BPM 2019 Blockchain and CEE Forum, Vienna, Austria, September 1–6, 2019, Proceedings 17*, pages 280–295. Springer, 2019.
- [20] Manling Li, Shiyu Zhao, Qineng Wang, Kangrui Wang, Yu Zhou, Sanjana Srivastava, Cem Gokmen, Tony Lee, Li Erran Li, Ruohan Zhang, Weiyu Liu, Percy Liang, Li Fei-Fei, Jiayuan Mao, and Jiajun Wu. Embodied agent interface: Benchmarking Ilms for embodied decision making. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 15, 2024, 2024.
- [21] Zelong Li, Shuyuan Xu, Kai Mei, Wenyue Hua, Balaji Rama, Om Raheja, Hao Wang, He Zhu, and Yongfeng Zhang. Autoflow: Automated workflow generation for large language model agents. *arXiv preprint arXiv:2407.12821*, 2024.
- [22] Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025.
- [23] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019.
- [24] Fionn Murtagh and Pedro Contreras. Algorithms for hierarchical clustering: an overview. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 2(1):86–97, 2012.
- [25] Frank Nielsen and Frank Nielsen. Hierarchical clustering. *Introduction to HPC with MPI for Data Science*, pages 195–211, 2016.
- [26] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. Chatdev: Communicative agents for software development. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 15174–15186. Association for Computational Linguistics, 2024.
- [27] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, 2016.
- [28] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *CoRR*, abs/2402.03300, 2024.
- [29] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023.

- [30] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Cote, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. {ALFW}orld: Aligning text and embodied environments for interactive learning. In *International Conference on Learning Representations*, 2021.
- [31] Olly Styles, Sam Miller, Patricio Cerda-Mardini, Tanaya Guha, Victor Sanchez, and Bertie Vidgen. Workbench: a benchmark dataset for agents in a realistic workplace setting. In *First Conference on Language Modeling*, 2024.
- [32] Harsh Trivedi, Tushar Khot, Mareike Hartmann, Ruskin Manku, Vinty Dong, Edward Li, Shashank Gupta, Ashish Sabharwal, and Niranjan Balasubramanian. Appworld: A controllable world of apps and people for benchmarking interactive coding agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16022–16076, 2024.
- [33] 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.
- [34] Zora Zhiruo Wang, Jiayuan Mao, Daniel Fried, and Graham Neubig. Agent workflow memory. *arXiv preprint arXiv:2409.07429*, 2024.
- [35] Judith Wewerka and Manfred Reichert. Robotic process automation—a systematic literature review and assessment framework. *arXiv preprint arXiv:2012.11951*, 2020.
- [36] Michael Wornow, Avanika Narayan, Krista Opsahl-Ong, Quinn McIntyre, Nigam H Shah, and Christopher Re. Automating the enterprise with foundation models. *arXiv preprint arXiv:2405.03710*, 2024.
- [37] Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. Travelplanner: a benchmark for real-world planning with language agents. In Proceedings of the 41st International Conference on Machine Learning, pages 54590–54613, 2024.
- [38] 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.
- [39] John Yang, Akshara Prabhakar, Karthik R Narasimhan, and Shunyu Yao. Intercode: Standardizing and benchmarking interactive coding with execution feedback. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.
- [40] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022, 2022.
- [41] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*.
- [42] Yining Ye, Xin Cong, Shizuo Tian, Jiannan Cao, Hao Wang, Yujia Qin, Yaxi Lu, Heyang Yu, Huadong Wang, Yankai Lin, Zhiyuan Liu, and Maosong Sun. Proagent: From robotic process automation to agentic process automation. *CoRR*, abs/2311.10751, 2023.
- [43] Zhen Zeng, William Watson, Nicole Cho, Saba Rahimi, Shayleen Reynolds, Tucker Balch, and Manuela Veloso. Flowmind: automatic workflow generation with llms. In *Proceedings of the Fourth ACM International Conference on AI in Finance*, pages 73–81, 2023.
- [44] Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. Expel: LLM agents are experiential learners. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan, editors, *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 19632–19642. AAAI Press, 2024.

- [45] Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar prompting with memory for computer control. *arXiv preprint arXiv:2306.07863*, 2023.
- [46] Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents. In *The Twelfth International Conference on Learning Representations*, 2024.

# **NeurIPS Paper Checklist**

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

# IMPORTANT, please:

- Delete this instruction block, but keep the section heading "NeurIPS Paper Checklist",
- · Keep the checklist subsection headings, questions/answers and guidelines below.
- Do not modify the questions and only use the provided macros for your answers.

#### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The main claims made in the abstract and introduction do accurately reflect our paper's contributions and scope.

# Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
  contributions made in the paper and important assumptions and limitations. A No or
  NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

# 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: See Appendix A.

#### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

# 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

12.1

Justification: Our paper does not contain theoretical results.

# Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

# 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: See Sections 4 and Appendix Hyperparameters.

# Guidelines:

• The answer NA means that the paper does not include experiments.

- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

# 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We include our source code and data in the supplemental material submission, and we outline the data generation procedure, the evaluation protocol, the training regime, and everything else necessary for reproduction either in the main body of the paper or in the appendix.

# Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be
  possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
  including code, unless this is central to the contribution (e.g., for a new open-source
  benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.

- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

#### 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See Section 4.1 and Appendix.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
  material.

# 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: The results we report are the average of multiple runs.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

#### 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: See Section 4.1.

# Guidelines:

• The answer NA means that the paper does not include experiments.

- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

#### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We have carefully checked the NeurIPS Code of Ethics and ensure our code aligns with the request.

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

# 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We provide the broader impacts of our work in Appendix A.

#### Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

# 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our paper does not have high risk PLMs or datasets.

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
  necessary safeguards to allow for controlled use of the model, for example by requiring
  that users adhere to usage guidelines or restrictions to access the model or implementing
  safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
  not require this, but we encourage authors to take this into account and make a best
  faith effort.

# 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We have cited all the corresponding papers.

#### Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
  package should be provided. For popular datasets, paperswithcode.com/datasets
  has curated licenses for some datasets. Their licensing guide can help determine the
  license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We have provided comprehensive and detailed documentation.

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

#### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: Our paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: Our paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

# 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [No]

Justification: In our research, the LLMs are used solely for writing, editing, and formatting purposes. It does not impact the core methodology, scientific rigor, or originality of the research. Therefore, we declare that the usage of the LLM does not affect the main contributions of the paper.

# Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

# **A** Limitation and Broader Impacts

**Limitation.** While the framework proposed in this paper represents a notable advancement in LLM agents, it also has certain limitations that warrant discussion. First, the proposed framework requires offline data in order to organize hierarchical MetaFlows. Second, the rejection sampling method and reinforcement learning approach employed in this paper rely on verifiable rewards, which limits their applicability in some open-ended evaluation tasks. Therefore, further exploration of these methods' applicability in such tasks is required.

**Broader Impact.** This paper introduces research aimed at enhancing LLM agent capabilities through experience reuse. From a societal impact perspective, while we have developed a generic LLM-based autonomous agent, the presence of biased offline datasets may lead to decisions with suboptimal outcomes. Additionally, there is potential for autonomous agents to be misused in malicious applications. In response to these concerns, all data used in this paper is publicly available and does not involve private information. Moreover, the proposed framework should not be used for any malicious purposes.

# **B** Reward Dynamics During RL Training

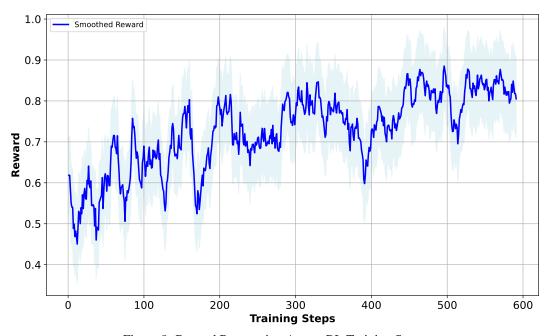


Figure 8: Reward Progression Across RL Training Steps.

As shown in Figure 8, despite variations in difficulty across different samples and the inherent uncertainty in the operation of the LLM agent leading to fluctuations in reward, the model's reward demonstrates a consistent improvement over the course of the three epochs in RL training.

Intuitively, during the RL process, the model generates various MetaFlows through rollouts. We calculate the reward by executing the MetaFlows at the leaf nodes, which is based on test case matching. This reward avoids the risk of reward hacking and is identical to the calculation process used during actual deployment. These rewards are then involved in the calculation of fitting weights for different trajectories. MetaFlows with higher rewards receive greater fitting weights, which steers the model's optimization towards these samples. On the other hand, MetaFlows with lower rewards are trained to avoid generating such samples.

For example, during the MetaFlowGen's rollout, two MetaFlows may be generated: one that does not meet the format specifications (resulting in a reward of -1) and another that adheres to the specifications, featuring more static steps and more detailed dynamic step instructions (resulting in a

# Algorithm 2 Deep Matching Traversal

```
Require: Hierarchical Experience Tree \mathcal{T}, task q, matching model \mathcal{P}(q|\mathcal{Q},\mathcal{M})
Ensure: Best matching MetaFlow instance
 1: function FindBestMatch(node, q, depth)
 2: if \mathcal{P}(q|\mathcal{Q}_{\text{node}}, \mathcal{M}_{\text{node}}) = \text{False then}
3:
        return None
4: end if
 5: best\_node \leftarrow node, max\_depth \leftarrow depth
 6: for each child c of node do
        match \leftarrow \text{FindBestMatch}(c, q, depth + 1)
7:
        if match \neq None then
 8:
9:
           child\_depth \leftarrow Depth \ of \ match
           if child\_depth > max\_depth then
10:
               best\_node \leftarrow match, max\_depth \leftarrow child\_depth
11:
12.
           end if
13:
        end if
14: end for
15: return best_node
16: Start from root: best\_node \leftarrow FindBestMatch(\mathcal{T}_{root}, q, 0)
17: return execute(q, \mathcal{M}_{best\_node})
```

reward of 1). The GRPO algorithm adjusts the probabilities accordingly—reducing the likelihood of generating non-compliant samples and increasing the probability of generating valid samples.

# C Deep Matching Traversal

Algorithm 2 illustrates the Deep Matching Traversal process, which recursively searches the hierarchical experience tree to find the deepest node whose experience best matches the query.

# **D** Dataset Statistics

For WorkBench, since the original repository only provided golden solutions without complete workflows, we first prompted Qwen-32B-Instruct and filtered responses that matched the golden solutions to obtain offline data. For AppWorld, we used the officially partitioned training set as the offline data. The dataset statistics are summarized in Table 5.

Metric	WorkBench	AppWorld
Offline Data Size	237	90
Test Data Size	353	57
SFT Data Size	1,102	1,092
RL Data Size	121	247

Table 5: Data statistics.

# E Implementation of Baseline Methods

The implementation details of the baseline methods are as follows.

**Trajectory-based baseline.** For each input task, we first encode its description using *all-MiniLM-L6-v2* to obtain a task embedding. Cosine similarity is then computed between this embedding and all task embeddings in the offline dataset. The most similar task is selected, and its description and solution trajectory are used as in-context learning (ICL) examples in the LLM agent's prompt.

**Guideline-based baseline.** We first prompt the same LLM to summarize guidelines from offline trajectories. For a given task, we again encode its description using *all-MiniLM-L6-v2* and find the most similar task via cosine similarity. The retrieved task's guideline and randomly selected ICL examples are then incorporated into the LLM agent's prompt.

The guideline summarization prompt is:

Please summarize the experience in natural language from the following examples so that you can use this experience in subsequent similar tasks. Please don't rehash the code or discuss code details. query: {query} solution: {trajectory}

# F Effect of Distance Metric on Experience Tree Quality

MetaFlowLLM is a general-purpose framework not tied to a specific distance metric. While our main experiments use task embeddings from *all-MiniLM-L6-v2* with cosine similarity, the framework also supports alternatives such as trajectory embeddings, Jaccard similarity, or rule-based measures.

As shown in Table 6, both trajectory- and task-based embeddings improve model performance over baselines. The slight gap between them likely arises from noise in trajectory representations caused by tool call responses.

Model	%Acc↑	%SE↓
Qwen2.5-7B	14.3	43.0
+ Trajectory Embedding	15.2	40.8
+ Task Embedding	17.2	46.8
Qwen2.5-32B	19.2	47.2
+ Trajectory Embedding	24.9	46.5
+ Task Embedding	28.3	40.4
GPT-4o-mini	19.6	43.0
+ Trajectory Embedding	25.1	37.1
+ Task Embedding	26.1	33.6

Table 6: Impact of Distance Metric in MetaFlowLLM.

# G Case Study of MetaFlows

In this section, we present a typical MetaFlow example on WorkBench.

```
"type": "dynamic",
    "instruction": "Extract the comparison operator, threshold,
        and time period from the user query and calculate the
       start and end date. The time period is natural language
        like '1 week' or '2 weeks'. Output the start and end
       dates as ISO 8601 format.",
   "outputs": [
      "time_min",
      "time_max",
      "comparison_operator",
      "threshold"
   ]
 },
   "type": "static",
   "action": "analytics.total_visits_count",
   "params": {
      "time_min": "${time_min}",
     "time_max": "${time_max}"
   },
    "output": "total_visits_data"
 },
```

```
{
  "type": "dynamic",
  "instruction": "Check if the total visits data meets the
    condition specified by the comparison operator and
    threshold. If so, generate a line chart using the
    analytics.create_plot function with the value_to_plot
    set to total_visits.",
  "outputs": [
       "plot_file_path"
  ]
}
```

# **H** Inference Time Comparison

Table 7 presents the inference time comparison between the standard ReAct agent and the proposed ReAct+MetaFlowLLM framework on the WorkBench benchmark. From the results, we observe that incorporating MetaFlowLLM incurs only a minor increase in total inference time (approximately 9%) compared to the vanilla ReAct agent, while achieving a substantial improvement in accuracy. This demonstrates that MetaFlowLLM offers a favorable trade-off between efficiency and performance, making it a practical enhancement for real-world multi-step reasoning systems.

Table 7: **Inference time comparison on WorkBench.** Both methods use Qwen2.5-32B-Instruct as the underlying model.

Method	Avg. Retrieval Time (s)	Avg. Running Time (s)	Accuracy (%)
ReAct	0.0	12.1	19.2
ReAct + MetaFlow	1.7	11.5	28.3

# I Comparison with AWM

We conduct experiments using application type as the categories in AWM [34], and the experimental results on AppWorld are presented in the Table 8. These results demonstrate that AWM indeed provides improvements over the trajectory-based and guideline-based baselines in this paper. Notably, MetaFlowLLM consistently outperforms AWM across all backbones.

Table 8: Performance comparison of various models and agent types on AppWorld (%).

Model	Agent Type	Configuration						
		Base	w/ Traj.	w/ Guideline	w/ MetaFlow	w/AWM		
Qwen2.5-7B	Reflexion ReAct	7.6 6.4	32.7 37.4	13.5 17.0	37.4 43.9	35.6 40.4		
0 25 220	Reflexion	12.9	33.9	22.8	43.9	25.7		
Qwen2.5-32B	ReAct	22.2	49.1	38.6	50.9	<u>50.3</u>		
GPT-4o-mini	Reflexion ReAct	13.5 7.0	37.4 22.2	20.5 7.0	45.6 35.7	37.4 22.8		

# J Prompt Design

The MetaFlow generation prompt was manually designed through iterative refinement during preliminary experiments, while agent prompts were adapted from the official prompts of their respective source benchmarks.

#### J.1 MetaFlow Generation Prompt

#### SYSTEM

You are tasked with analyzing two provided examples, each consisting of a user query and its corresponding workflow. Based on these examples, your objective is to generalize a meta workflow that can solve similar problems using a combination of static executable code and dynamic LLM agent configurations. The meta workflow you generalize should not only summarize the given samples, but it should be as general as possible so that it can solve similar problems by reusing static code and utilizing dynamic agents.

The meta workflow runs sequentially. When encountering static code, it is executed directly. When encountering a dynamic agent's JSON configuration file, an LLM agent is triggered to generate Python code to solve the specific task. The Python code generated by the LLM agent and the static code will operate within the same compilation environment, allowing them to reference and interact with each other.

# # Key Requirements:

# ## Static Code:

The static code segments must be fully executable and should not contain any placeholders or dynamic interpolation logic (e.g., query-specific formatted strings).

These segments should perform reusable tasks that are common across similar problems.

Considering that the examples provided may vary significantly, making it impossible to reuse any static code, static code can be omitted. However, please avoid generating static code filled with comments, as this lacks practical significance. The explanatory text can be placed entirely in the configuration file of the dynamic agent.

#### ## Dynamic Configurations:

The dynamic agent should complete task-specific tasks and be able to link the static code before and after, ensuring that they run without errors. The dynamic portions of the meta workflow should be defined in a JSON configuration file.

The JSON configuration must include:

- The `task\_description` that clearly defines the queryspecific logic using placeholders (e.g., {criteria}) to represent components of the user's query.
- The `expected\_final\_state` upon task completion, such as how variables are updated to enable integration with subsequent static code. For example, if a later static Python code requires variables to be defined during a dynamic agent's execution, this should be clearly indicated.

#### ## Meta Query:

The meta query should describe the type of problem that the meta workflow can solve, without including all specific details from the input queries. It should use placeholders (e.g., [criteria]) to generalize the problem.

# ## Meta Workflow Structure:

The meta workflow should alternate between static code blocks and dynamic JSON configurations.

- Static code blocks must be self-contained and executable without modification. Code consisting entirely of comments is not acceptable.
- Dynamic JSON configurations should only contain placeholders for query-specific logic and should not interfere with the execution of static code.

```
USER:
**Now Summary the Following Examples:**
**API Documentation**:
{api_docs}
**Sample1**:
**Query**:
{query1}
**Workflow**:
{workflow1}
**Sample2**:
**Query**:
{query2}
**Workflow**:
{workflow2}
Output:
**Meta Query**:
[Provide the generalized meta query here]
**Meta Workflow**:
[Provide the generalized meta workflow, which should include
   the static code enclosed appropriately and a dynamic agent
   configuration in JSON format. The configuration should be
   clearly separated by `---` and enclosed within triple
   backticks (```) to maintain structure and clarity.]
```

# J.2 Agent Prompt Design in WorkBench

#### SYSTEM:

You are a step-by-step task-solving assistant. At each stage, you are only allowed to complete the current subtask based on the available variables, tools and instructions.

Today's date is Thursday, November 30, 2023 and the current time is 00:00.

Please remember the current date and time. Please use today's date for all dates unless otherwise specified. Meetings must not start before 9am or end after 6pm.

Please respond helpfully and accurately to the user.

All the information is in the given query or you can get it using the given tools. Please do not assume or will ask for my help.

#### USER:

You are solving a multi-step query within a pre-defined workflow. The full user query is:

{total\_task}

And the pre-defined workflow is:
{workflow}

But you should now focus \*\*only\*\* on the following subtask: {step['instruction']}

Here are the historical Logs prior to the completion of this sub\_task:

{history\_logs}

Use the current variables: {context}

Do not perform any part of the full task beyond this subtask. Do not call other tools or add outputs unless explicitly instructed.

#### J.3 Agent Prompt in AppWorld.

#### SYSTEM:

You are a super intelligent AI Assistant whose job is to complete day-to-day tasks by writing code to interact with apps on behalf of your supervisor. Use API documentation to understand how to interact with the apps.

You will undertake a \*multi-step conversation\* using a python REPL environment. That is, you will write the python code and the environment will execute it and show you the result, based on which, you will write python code for the next step and so on, until you've achieved the goal. This environment will let you interact with app/s using their associated APIs on my behalf.

Currently, you need to address the final sub-task of a total task while ensuring the variables used are consistent with previous interactions. We have executed some history code in previous steps; that history code is provided here, and you may assume that any relevant state or variables from it are in memory. Proceed under the assumption that you can use or refer to prior variables as needed.

#### USER:

# App-wise API Documentation:
 ``yaml
{api\_documentation\_string}

# ASSISTANT:

Understood.

```
USER:
You have access to the following imports already available in
your coding environment.
{available_imports}
These APIs should be called as python functions through the `
   apis` object. E.g., `apis.supervisor.show_profile(...)` for
   the `show_profile` API of the `supervisor` app. Note that
   you already have the `apis` object available in the
   execution environment, so do NOT to create it yourself.
You can also import anything from the Python *standard* library
    via code. However, you cannot use any systems-related
   packages, like os, shutils, etc, as they are disabled.
Let's say you want to obtain supervisors' app account passwords
    and get login access_token for one of their apps, let's
   say Spotify. You can write the following snippet as part of
    your code:
```python
# I should use supervisor.show_profile to get the supervisor's
   account name and password,
# then pass it to spotify.login API in terms of username and
   password
supervisor_profile = apis.supervisor.show_profile()
supervisor_passwords = __CURLY_OPEN__
    account_password["account_name"]: account_password["
       password"]
    for account_password in apis.supervisor.
       show_account_passwords()
__CURLY_CLOSE__
spotify_access_token = apis.spotify.login(
    username=supervisor_profile["email"],
    password=supervisor_passwords["spotify"],
)["access_token"]
# ... remaining code uses spotify_access_token variable as
required.
Notice how the arguments passed to the APIs and outputs parsed
   from their outputs are as per the API documentation.
After you have completed the total task, you must call `apis.
   supervisor.complete_task`.
If the task is a question (e.g., "How many songs are in the
   Spotify queue?"), it must be called with an `answer`
   argument with an appropriate value. Use words or numbers
   only as answers, not full sentences, e.g., "10" in this
   case.
If the task is not a question, "Start my Spotify music player
   .", the `answer` argument should not be passed, or its
   values should be `None`.
```

ASSISTANT: Got it.

# # History Code and Feedback {history} # Total Instruction {task\_instruction} # Subtask Instruction {subtask\_instruction} Write the code to complete this subtask. Only generate valid Python code. It must be within markdown-styled ("``") code blocks. Do NOT say or explain ANYTHING else. # APIs allowed to Use {{required\_apis | join(", ")}}

# Remember you:

- must only use APIs from the above list passing arguments and parsing outputs as per the provided documentation.
- must make all decisions autonomously on your own, and not wait or ask for anything external.
- must call `apis.supervisor.complete\_task` at the end as per the above instruction.
- do not have access to any other packages except what is provided above and is part of the Python standard library.
- You shouldn't write all the code at once. Write small chunks of code and only one chunk of code in every step to make full use of the environment feedback. Make sure everything is working correctly before making any irreversible change.
- Before each time you write a specific code, you should think about the current situation in the form of a code comment and put it in the code blocks.