

# IMPROVING LARGE LANGUAGE MODEL BASED MULTI-AGENT FRAMEWORK THROUGH DYNAMIC WORKFLOW UPDATE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Multi-agent frameworks powered by large language models (LLMs) have demonstrated great success in automated planning and task execution. However, the effective adjustment of workflows during execution has not been well-studied. A flexible workflow is crucial, as in many real-world scenarios, the initial plan must adjust to unforeseen challenges and changing conditions in real-time to ensure the efficient execution of complex tasks. In this paper, we define *workflows activity-on-vertex (AOV) graphs*. We continuously refine the workflow by dynamically adjusting task allocations based on historical performance and previous AOV with LLM agents. To further enhance system performance, we emphasize *modularity* in workflow design based on measuring parallelism and dependence complexity. Our proposed multi-agent framework achieved efficient sub-task concurrent execution, goal achievement, and error tolerance. Empirical results across different practical tasks demonstrate dramatic improvements in the efficiency of multi-agent frameworks through dynamic workflow updating and modularization.

## 1 INTRODUCTION

Large Language Models (LLMs) (Significant Gravitas; Zhou et al., 2023) show remarkable abilities to understand and generate human-like text. Recent advances have significantly enhanced their capability to emulate human reasoning (Sun et al., 2024), indicating a promising future for LLM-based reasoning. With the powerful ability to deal with a variety of natural language processing tasks, these models underpin a wide range of applications, from conversational agents (Ye et al., 2024) and content creation tools (Yao et al., 2023) to advanced analytics and decision-making systems (Ramesh et al., 2021; Wang et al., 2023). Building upon this foundation, a key advancement is the development of *multi-agent systems empowered by LLMs* (Liu et al., 2023; Li et al., 2023; Hong et al., 2024; Wu et al., 2024; Wang et al., 2024; Chen et al., 2024) where multiple LLM-based agents collaborate to address the same task, leveraging their collective reasoning and planning abilities to automate and optimize task execution processes.

Existing LLMs-based multi-agent systems define LLM as an agent and agents are collaborated with each others via manually designed or LLM-generated prompts. Specifically, MetaGPT (Hong et al., 2024) focused on programming tasks by leveraging Standardized Operating Procedures (SOPs) (Wooldridge & Jennings, 1998; DeMarco & Lister, 2013; Belbin, 2010). It predefined distinct roles such as product manager, project manager, and engineer. For each role, an LLM agent is initialized, and these agents follow a strict and sequential workflow to execute sub-tasks. CAMEL (Li et al., 2023) is designed to complete a variety of tasks. It requires users to predefined two agents. These agents interact and execute tasks sequentially, with each agent taking on specific responsibilities. AutoGen (Wu et al., 2024) is also aimed at completing diverse tasks. Unlike CAMEL, AutoGen can automatically create an agent list with different roles based on the task requirements. These agents execute tasks sequentially following the order in the list.

Building upon the strengths of current multi-agent systems, our work aims to further improve existing general-purpose multi-agent systems by enabling *dynamic updating workflows* during task execution and encouraging *modularity* in workflows when planning the workflows.



Figure 1: Comparative evaluations among four frameworks—AutoGen, CAMEL, MetaGPT, and Flow (ours)—across two tasks, present notable differences in performance. For the left task, AutoGen, CAMEL, and MetaGPT only managed to produce basic designs lacking in completeness while Flow excelled by creating a fully developed and well-structured website. For the right task, our Flow demonstrated superior capability by successfully generating a working game with a clear and intuitive interface while the other frameworks struggled to deliver fully functional or correct code.

Specifically, *dynamic updating workflow* allows to adjust *sub-task allocations* and *agent roles* in real-time based on ongoing performance feedback and changing conditions. This capability *ensures* that the system remains responsive and efficient even when faced with unexpected obstacles. For instance, if an agent encounters a roadblock in data preprocessing, dynamic updating allows the system to reassign the sub-task to another agent or introduce a new sub-task to overcome the challenge. This adaptability is essential for maintaining robustness and ensuring the seamless execution of complex tasks.

**Modularization** in system design involves dividing a system into separate, independently operating modules, each responsible for specific functionalities (Baldwin & Clark, 1999). A highly modularized system allows each module to be developed, managed, and executed in isolation, which simplifies system design and enhances adaptability. In our context, modularization refers to the decomposition of a complex task into smaller, interchangeable sub-task modules. A highly modularized workflow enables sub-tasks to execute concurrently, without bottlenecks from other parts of the workflow. It directly improves operational efficiency of multi-agent frameworks. In addition, modularity dramatically enhances the ease of dynamic updating. When workflows are highly modularized, the dependency complexity between sub-tasks is small. Therefore, updating one sub-task does not necessitate changes in others, allowing for a small adjustments. For instance, if an agent responsible for data preprocessing encounters an unexpected obstacle, the system can dynamically introduce a new sub-task to address the issue with a little influence to the rest of the workflow.

In this paper, we have improved existing multi-agent systems by fulfilling modularity and dynamic updating workflow. Our system allows agents to run their sub-tasks in a parallel manner while enabling effective dynamic updates to workflows simultaneously by formulate the entire workflow as an Activity-on-Vertex (AOV) graph, which is a directed acyclic graph (DAG) where each sub-task is represented as a node with its status and generated logs, and the directed edges capture dependencies between sub-tasks. To encourage a modularized workflow design from the beginning, we generate multiple candidate AOV graphs for the task. These candidates are then evaluated based on their degree of parallelism and the complexity of their dependencies. The AOV graph with The highest parallelism and lowest dependency complexity will be selected.

During task execution, our system dynamically updates the workflow when a sub-task is found failure (more detail on Fig. 2: *Running and Tracking status*). Updating the system involves modifications to task allocations and agent roles based on ongoing performance data and current workflows. As our AOV-based workflow are encouraged to have high modularity, updating one module does not necessitate changes in others, allowing for localized adjustments during workflow updates (more detail on Fig. 2: *Refining*). Similar to the initial workflow generation, multiple AOV graphs are generated and the one with the *highest parallelism* and *lowest dependency complexity* is selected during the dynamic updates. This iterative workflow refinement process ensures a good capability of adapting to new challenges and evolving objectives throughout task execution without compromising overall performance. Our key contributions are as follows:

- 1). We identify a key limitation in current multi-agent frameworks—the reliance on rigid, sequential workflows, which hinder adaptability and scalability. We emphasize the importance of workflow modularity and dynamic updates, facilitated by global information, to enhance system robustness and flexibility.
- 2). We introduce a practical multi-agent framework designed to support concurrent sub-task execution and dynamic workflow updates. This approach allows agents to adapt efficiently to unexpected challenges, thus improving system scalability and facilitating easier maintenance.
- 3). Our experiments demonstrate that the proposed framework dramatically outperforms existing frameworks. By handling complex tasks with greater resilience and quality, our framework leads to notable improvements in performance metrics.

## 2 RELATED WORK

**LLM based Task Decision-Making** Recent developments in LLM-driven task decision-making have focused on enhancing the reasoning and planning abilities of agents. previous approach like ReAct (Yao et al., 2023) which iteratively generates thoughts and actions based on current observations until task completion. This framework integrates action-taking with reasoning, allowing agents to perform complex tasks in dynamic environments. Reflexion (Shinn et al., 2023) further improves this by incorporating self-reflection, where the agent evaluates and adjusts its reasoning during execution. ADAPT (Prasad et al., 2023) introduces recursive task decomposition, enabling LLM-based agents to break tasks into smaller subtasks, leading to improved task execution flexibility. However, these approaches often overlook dynamic task reallocation, particularly in multi-agent settings, which is where our work extends the current research.

**LLM based Multi-Agent Frameworks** Multi-agent frameworks have long been employed for task execution in distributed environments, with recent advancements leveraging LLMs to improve coordination and decision-making. Current frameworks like MetaGPT (Hong et al., 2024) and CAMEL (Li et al., 2023) employ structured workflows where multiple agents collaborate to accomplish complex tasks. However, these frameworks typically rely on static workflows, which limit their ability to dynamically adapt to changes in the task environment. Recent works like AutoGen (Wu et al., 2024) attempt to address this by introducing more flexible agent collaboration mechanisms. Our approach takes this further by introducing AOV-based dynamic task planning, enabling more efficient subtask allocation and independent task execution, thus improving adaptability in complex scenarios.

## 3 METHOD

Our proposed framework enhances multi-agent frameworks powered by LLMs by introducing modularity and dynamic workflow updating. This section details how we achieve these features.

**Formulating a Workflow as an AOV Graph** Activity on Vertex (AOV) graph is a type of directed acyclic graph (DAG) where vertices represent tasks and edges denote precedence relations (Bondy & Murty, 2011). AOV graphs are crucial in project scheduling and management (Moder et al., 1983; Taha, 2017), helping planners visualize dependencies and sequence tasks efficiently.

Inspired by that, we define Multi-Agent workflow as an AOV graph where vertices represent sub-tasks, with its edges denoting dependencies between these sub-tasks. Let  $G = (V, E, A)$  denote the AOV Graph, where  $V$  is the set containing all sub-tasks (nodes),  $E \subseteq V \times V$  represents the set of directed edges indicating sub-task dependencies, and  $A$  represents a set of agents for all sub-tasks. Each agent  $a_j \in A$  is associated with a role  $s_j$  and is responsible for executing a subset of tasks  $\mathcal{T}_j \subseteq V$ . We also generate sub-tasks and each directed edge  $e_{ij} = (v_i, v_j) \in E$  indicates that sub-task  $v_i$  must be completed before sub-task  $v_j$  can be started.

Note that AutoGen (Wu et al., 2024) also automatically generates sub-tasks and agents. However, the sub-tasks are designed to be executed *sequentially*. For Flow, we allow for the generation of complementary sub-tasks that can run *in parallel*. This distinction enhances our system’s ability to handle multiple tasks simultaneously, which reduces overall process time and increase efficiency.

**Modularity in a Workflow** Modularity in system design (Baldwin & Clark, 1999) involves dividing a system into separate, independently operating modules, each responsible for specific functionalities, allowing focus on individual components without affecting the entire system. In the context of workflows, we advocate for the creation of sub-tasks that can be executed independently. Modularity is essential for scalability and flexibility in workflows. By reducing dependency complexity, the system can more easily adapt to changes, such as the introduction of new tasks or the reassignment of existing ones, without requiring extensive restructuring.

To encourage modularity in the generated AOV Graph, we define two quantitative measures that evaluate parallelism and dependency complexity, respectively. Parallelism measures the extent to which tasks can be executed concurrently. Let  $S_t$  represent the set of tasks executed at step  $t$ . Let  $T$  be the total number of steps (the maximum depth of The DAG). Given a AOV Graph  $G = (V, E, A)$ , the degree of parallelism at a specific step  $t$  is defined as the average ratio of the number of tasks executed in that step to the total number of tasks:

$$P_{\text{avg}} = \frac{1}{T} \sum_{t=1}^T P(t), \text{ where } P(t) = \frac{|S_t|}{|V|}.$$

While  $P_{\text{avg}}$  provides a measure of parallelism, it is insufficient to fully capture the modularity that arises when sub-tasks can be executed independently. Consider two workflows, both containing the same sub-tasks  $\{A, B, C, D\}$ . For Workflow 1, the task dependencies are defined as:  $A \rightarrow C, B \rightarrow C, A \rightarrow D, B \rightarrow D, C \rightarrow D$ . In contrast, Workflow 2 has dependencies:  $A \rightarrow C, B \rightarrow C, C \rightarrow D$ . Although both workflows exhibit the same level of parallelism, Workflow 2 is structurally simpler in terms of task dependencies, as it contains fewer edges.

To account for this complexity, we measure the dependency structure by analyzing the degree distribution within the task graph. For each task  $v_i$ , the degree  $\deg(v_i)$  reflects the number of direct connections it has in the graph  $G$ . The average degree  $\bar{d}$  is computed as:

$$\bar{d} = \frac{1}{|V|} \sum_{v_i \in V} \deg(v_i),$$

where  $|V|$  is the number of tasks (vertices) in the graph. The complexity of task dependencies is then quantified by the standard deviation of the degree distribution:

$$C_{\text{dependency}} = \left( \frac{1}{|V|} \sum_{v_i \in V} (\deg(v_i) - \bar{d})^2 \right)^{\frac{1}{2}}.$$

This measure reflects the variability in the number of dependencies each task has, providing insight into the overall complexity of the workflow structure.

Task dependencies alone are insufficient to fully capture the modularity that allows sub-tasks to be executed independently. Consider Workflow 3:  $A \rightarrow B \rightarrow C \rightarrow D$ , which may have a similar dependency complexity to Workflow 2. However, Workflow 2 provides greater modularity and separation of tasks, highlighting the importance of evaluating both dependency complexity and modularity to fully assess and promote effective workflow designs. Both measures are essential for ensuring that tasks can be executed in parallel while maintaining a well-structured, modular approach.

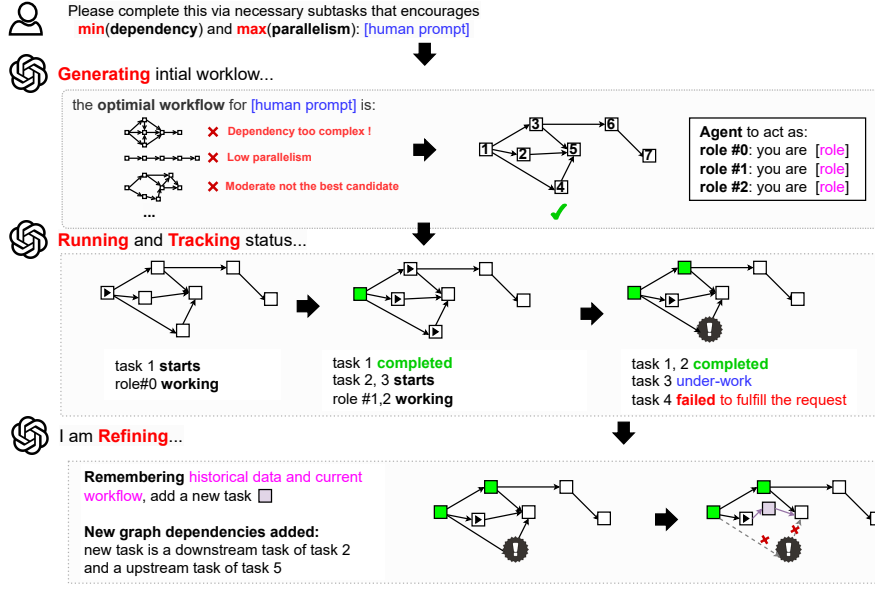
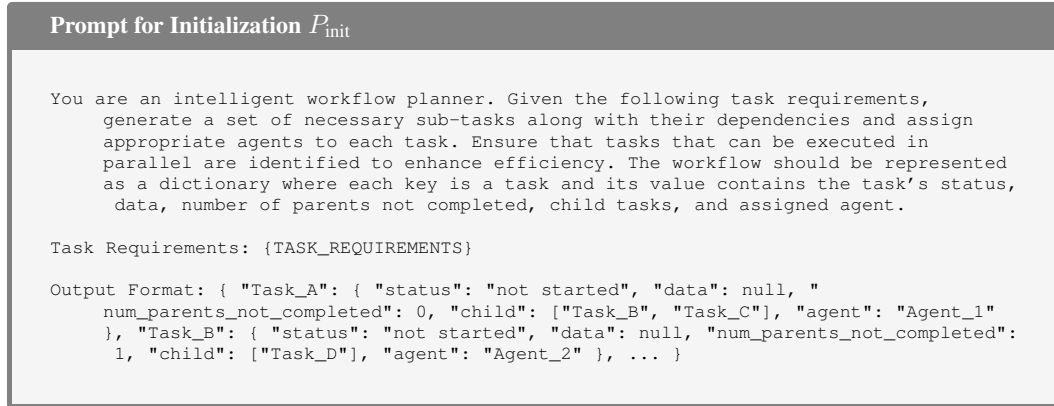


Figure 2: Our system encourages the modularity of the sub-tasks and allows agents to run their tasks in a parallel manner while enabling dynamic updates to workflows simultaneously.



**Generate an Initial AOV Graph** Given a task requirement  $T$ , firstly, we prompt a LLM  $f$  to generate a set of candidate AOV Graphs  $\{G_1, G_2, \dots, G_K\}$  based on the task requirements and our **Prompt for initialization**  $P_{init}$ , i.e.,  $\{G_1, G_2, \dots, G_K\} = f(P_{init}, T)$ . Each candidate AOV Graph  $G_k = (V_k, E_k, A_k)$  is evaluated using the measures of parallelism and dependency complexity. We prioritize the workflow with the highest parallelism score. If after the selection, the graph is not unique, we further select the one with the lowest dependency complexity.

**Execution Plan Generation and Agent Allocation** After we get the best candidate for the AOV graph, We begin by performing a topological sort on the task dependency graph to determine the order of task execution. The topological sort produces a linear ordering of the tasks  $\sigma : V \rightarrow \{1, 2, \dots, |V|\}$  such that for any edge  $(v_i, v_j) \in E$ ,  $\sigma(v_i) < \sigma(v_j)$ . The result is a sequence of task steps, where each step consists of tasks that can be executed in parallel. This execution plan minimizes the number of steps needed to perform while ensuring that all tasks are completed in the shortest possible time, adhering to their dependencies.

Each agent  $a_j \in A$  is associated with a set of sub-tasks  $\mathcal{T}_j \subseteq V$ , indicating the tasks that the agent is capable of handling. In our framework, we allow for the reuse of agents across different tasks based on their roles and time availability. However, if two sub-tasks  $v_p$  and  $v_q$  require the same agent  $a_j$

at the same step  $s_i$ , we create a clone of the agent, denoted  $a'_j$ , to run both sub-tasks simultaneously without increasing the wall time.

#### Prompt for Update $P_{update}$

You are an intelligent workflow updater. Based on the current workflow and the latest task progress data, update the workflow by adding, removing, or modifying sub-tasks as necessary. Ensure that the updated workflow maintains modularity and maximizes parallel execution.

Output Format: { "Task\_A": { "status": "not started", "data": null, ... } }

**Workflow Refinement and Dynamic Updating** Our dynamic updates are designed to be flexible, allowing modifications to task allocations including deletion, addition, editing, rerunning, and reassignment of agents. Without a modularity constraint, such flexibility would be difficult to implement. For instance, subtask dependencies can be very complex, and dynamically changing a task could necessitate redoing many existing tasks or incorporating many new tasks. With modularity, efficiency in our dynamic updating process is dramatically enhanced. Intuitively, when workflows are modular, updating one module does not necessitate changes in others.

We leverage Large Language Models (LLMs) as a global inspector to update an AOV graph based on global information. Specifically, given task requirements  $T$ , a prompt for update  $P_{update}$ , the current AOV Graph  $G^t$ , and generated data  $D^t$  containing the status of subtasks and the output of agents for running subtasks, the LLM continuously monitors task progress and dynamically modifies the graph when necessary. Similar to the initialization process, we also generate  $K$  candidate graphs:  $\{G_1^{t+1}, G_2^{t+1}, \dots, G_K^{t+1}\} = f(P_{update}, T, D^t)$ . We follow the same selection strategy as in initialization which prioritizes the workflow with the highest parallelism score. If after the selection, the graph is not unique, we further select the one with the lowest dependency complexity.

Note that with sufficient data and computational resources, we could further enhance our framework by fine-tuning LLMs with reinforcement learning (RL) for workflow generation. For example, the LLM would be trained to maximize a reward function designed around key performance indicators such as task completion speed, resource utilization, and minimization of workflow disruptions.

**Implementation** Our framework employs a dictionary-based structure,  $\tilde{G}$ , to efficiently manage and dynamically update workflows within a multi-agent system. This approach represents each task  $v$  in the workflow as a key in  $\tilde{G}$ , with the value being another dictionary that encapsulates various attributes of the task. The structure is specifically defined as:

$$\tilde{G}[v] = \{\text{"sub-task requirement"}, \text{"status"}, \text{"data"}, \text{"num\_parents\_not\_completed"}, \text{"child"}, \text{"agent"}\}.$$

Each task's dictionary includes attributes such as the sub-task requirement, current status (e.g., "not started", "in progress", "completed"), data relevant to the task, a count of uncompleted parent tasks to manage dependencies, a list of child tasks that depend on the current task's completion, and the agent assigned to the task. The choice of a dictionary-based structure for our workflow system is driven by its inherent simplicity and flexibility. This structure can be converted directly to JSON, and the organized information is easily readable and summarizable by large language models (LLMs).

Each task's execution readiness is determined by the attribute "num\_parents\_not\_completed". Tasks with a count of zero are eligible to run concurrently, leveraging our system's capability to handle parallel task execution effectively. Upon the completion of any task, we perform a systematic review to determine if the workflow requires refinement, ensuring that all dependencies are accurately accounted for and that the workflow remains aligned with project goals. Additionally, we do not rely solely on the status and "num\_parents\_not\_completed" counts reported by agents. These are always double-checked to prevent errors that could arise from inaccurate reporting by agents or unforeseen system anomalies. This rigorous verification process enhances the reliability and integrity of our workflow management system.



## 4 EXPERIMENTS

**Baselines** In all experiments, we compare our method, Flow, to the exists multi-agent frameworks *i.e.* (1) AutoGen (Wu et al., 2024), (2) Camel (Li et al., 2023), and (3) MetaGPT (Hong et al., 2024). In all our experiments, we use agents empowered by GPT-4o-mini (OpenAI, 2024) in a zero-shot setting.

**Experiment Design** We designed three diverse and engaging tasks to evaluate multi-agent collaboration frameworks: 1) website development, 2) LaTeX Beamer slide creation, and 3) interactive game development. The rationale behind selecting coding-based experiments is twofold. First, most multi-agent frameworks tend to favour coding and writing abilities, like MetaGPT (Hong et al., 2024). Using non-coding tasks may introduce bias. Second, coding tasks effectively demonstrate the system’s ability to assign agents and manage task allocation.

- *Development of a Gobang Game with Naive AI*: This task requires creating a Gobang (Five in a Row) game with a user interface and a simple AI opponent. Players can choose between black or white stones, with the UI clearly indicating turns and announcing the winner or draw when the game ends. This task demonstrates the system’s ability to handle modular design and task parallelism, as it involves coordinating game logic, AI implementation, and user interface development simultaneously.
- *Machine Learning Course Lecture Slides*: This task focuses on generating LaTeX slides covering reinforcement learning algorithms, including motivations, problem statements, intuitive solutions, and detailed mathematical equations. A specific page requirement is to test the system’s ability to follow instructions precisely. The task highlights the system’s parallel processing capabilities of simultaneous generation of content, formatting, and presentation structure. The structured format of LaTeX also tests how effectively the system manages modularity and concurrent tasks.
- *Development of a Comprehensive Website for ICLR 2025*: This task involves building a professional website for the International Conference on Learning Representations, hypothetically scheduled for San Francisco from April 27 to May 1, 2025. The website must feature key elements such as a detailed conference schedule and venue information with an interactive map. This task assesses each system’s ability to manage parallel workflows and modular components, including user interface design, functionality, and adherence to design guidelines, showcasing how well the system handles task decomposition and execution.

### 4.1 EVALUATIONS OVER THREE DESIGNED TASKS

**Evaluation Metrics** To conduct both quantitative and qualitative evaluations, we employed two metrics: *Success Rate* and *Human Rating*. ***Success Rate***: The Success Rate is a quantitative measure that ranges from 0 to 1. Assesses whether the multi-agent system successfully generates executable outputs that fully meet the task requirements. A higher score indicates a greater level of success in accurately fulfilling the task objectives. **For different tasks, it may have different evaluation metrics. The detailed description for each evaluation metric is defined in Appendix B.1, B.2 and B.3.**

***Human Rating***: Human ratings are used to evaluate the quality of the generated results in alignment with the task description. We gathered 50 participants with programming and machine learning backgrounds to rank the outcomes produced by different methods. the detailed description for how we take score is shown in the Appendix A

**Summary** We here give a summary of the performance of different methods over three tasks from Table 1, 2 and 3, comparing the overall score regarding the success rate and human rating. The overall score of Flow and human rating over three tasks, are (100, 4) on game design, (100, 3.33) on LaTeX writing, and (80, 3.28) on website design. Therefore, the average performance of Flow is 93% success rate and 3.54 over 4 satisfaction. Similarly, we have the average performance of AutoGen as (66.7, 2.75), MetaGPT as (71, 1.60), and CAMEL as (48.67, 2.12). Overall, our method Flow has finished tasks with the most satisfaction and the highest success rate. Information about Flow’s workflow on those task is in Appendix C

Table 1: Comparison of different LLM-based Multi-Agent frameworks on Gobang Game

Model	Success Rate (%)				Human Rating (1-4)
	Compilable	Intractable	Game Rule	Overall Score	
AutoGen (Wu et al., 2024)	80	60	40	60	2.26
MetaGPT (Hong et al., 2024)	100	100	20	73	1.24
CAMEL (Li et al., 2023)	40	40	0	27	2.50
<b>Ours</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>4.00</b>

#### 4.1.1 RESULT FOR GOBANG GAME

The experimental setup is thoroughly detailed in Appendix B.2 and the visualisation result is in Fig. 1. As shown in Table 1, our method gets 100 for all the aspects regarding success rate as well as the highest satisfaction from humans. More explanations for each method are as follows:

*AutoGen*: With the five tests, one trail failed to generate a valid output. Of the four successful attempts, one contained a code error that hindered normal execution, while another exhibited a bug in the game interface. The remaining two tests were completed successfully, though the chess pieces were displayed as the text “black” and “white”.

*MetaGPT*: After running MetaGPT five times, all attempts were successful and intractable. However, in four cases, a Tic-Tac-Toe game was generated instead of Gobang; out of these, the left one were functional, allowing both the user and AI to make moves and correctly terminate.

*CAMEL*: In all five trials, CAMEL was only successful twice. In the other attempts, the generated Python code was not executable. In the two successful trials, CAMEL successfully implemented correct termination conditions but had no AI component and terminated message.

*Flow*: After five rounds of testing, our system consistently generated successful outputs without any errors. The game functioned as expected, allowing both the player and the naive AI to take turns seamlessly. The game also ended correctly when either the board was fully occupied or one side achieved victory. In the game interface, the chess pieces were represented by actual black and white pieces, rather than text labels.

#### 4.1.2 RESULT FOR LATEX BEAMER WRITING

The experimental setup is thoroughly detailed in Appendix B.1 and the visualisation is in Fig. 10. Details of experimental results are presented in Table 5 with explanations as follows:

*AutoGen*: After five tests, AutoGen successfully generated outputs every time. However, one output failed to compile in LaTeX due to syntax errors, and in two instances, the outputs did not meet the required length. The remaining outputs met both the length and content requirements.

*MetaGPT*: In five trials, four of them successfully generated a valid LaTeX version, with the only error being related to writing Python code within the '.tex' file. In these four successful trials, all documents met the required content specifications, but the total page count fell short of the requirements of 30 pages or 20 pages.

*CAMEL*: Successfully generated five different '.tex' files that are valid and could be rendered into Beamer format. Each presentation contained the required information, including sections like motivation. However, none of them met the page count requirement of 30 pages or 20 pages.

*Flow*: After five tests, our system successfully generated outputs each time, and all outputs were able to compile in LaTeX. However, one output contained repetitive content. In the remaining valid outputs, the length of the Beamer presentations met the specified requirements, and all the content mentioned in the requirements was adequately covered.

#### 4.1.3 RESULT FOR WEBSITE DESIGN

Similar to the previous two, the detailed experiment set up is in Appendix B.3. We here illustrate the results in Table 3 as follows:



Table 2: Comparison of different LLM-based Multi-Agent frameworks on LaTeX Beamer writing

Model	Success Rate (%)				Human Rating (1-4)
	Compilable	Completeness	Page Limit	Overall Score	
AutoGen (Wu et al., 2024)	80	80	40	67	3.00
MetaGPT (Hong et al., 2024)	80	80	20	60	1.83
CAMEL (Li et al., 2023)	<b>100</b>	<b>100</b>	0	66	1.83
<b>Ours</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>3.33</b>

Table 3: Comparison of different LLM-based multi-agent frameworks on Website Design.

Model	Success Rate (%)				Human Rating (1-4)
	Compilable	Basic Information	Sections	Overall Score	
AutoGen (Wu et al., 2024)	80	80	60	73	2.62
MetaGPT (Hong et al., 2024)	<b>100</b>	<b>100</b>	40	<b>80</b>	1.72
CAMEL (Li et al., 2023)	80	80	0	53	2.02
<b>Ours</b>	80	80	<b>80</b>	<b>80</b>	<b>3.28</b>

*AutoGen*: AutoGen produced HTML and CSS files with key information displayed but lacks of details. Each section of the website contains only one or two sentences, lacking interactive functionality and necessary elements such as maps or tables.

*MetaGPT*: MetaGPT managed to create complete HTML and CSS, meeting basic functionality requirements and showcasing its code generation capabilities. However, the outputs were overly simplistic, missing significant content and key functional modules like the required venue and map.

*CAMEL*: CAMEL’s Outputs were executable in four out of five runs, though they did not include all the necessary elements, achieving all basic functions only. The system limits the communication can be only between two agents regardless of task complexity hindering its ability to fully complete complex website development tasks that require multi-task collaboration. Notably, one run generated complete HTML code but omitted the CSS file, preventing proper rendering of the website.

*Flow*: Flow gets 100% success rate. Furthermore, each section of the website featured detailed introductions and necessary interactive functionalities. For example, the venue section included travel information and local transportation options like airport, and accurately presented the conference location on a map. The registration section was fully functional, with complete table, input boxes, and a submission button.

## 5 WORKFLOW UPDATE

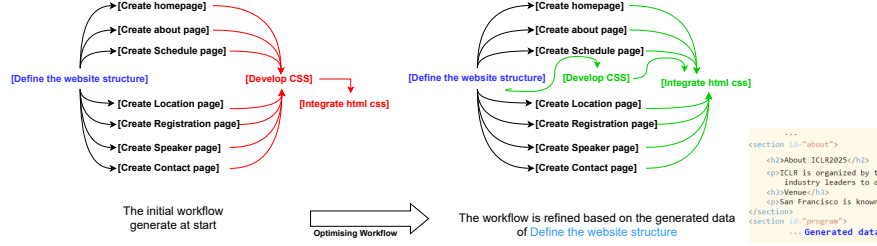
**Update based On Generated Data** Fig. 3(a) demonstrates the update process of Flow in the conference website creation example. Upon completing the first subtask, the system identifies potential changes and redundancies, triggering a restructuring process to enhance efficiency. Once the task "Define the website structure" is completed, the generated data, which includes HTML structures and elements is sufficient to proceed with the CSS creation. As a result, the workflow is updated to incorporate the development of CSS based on the completed "Define the website structure" task.

**Error handling** To evaluate the effectiveness of our updating mechanism, we intentionally introduced random masking to certain task outputs, replacing them with "none" before passing them to the next agent. We conducted five trials and recorded the success scores. Since other frameworks employ a sequential workflow, we limit the comparison to our own approach in this context.

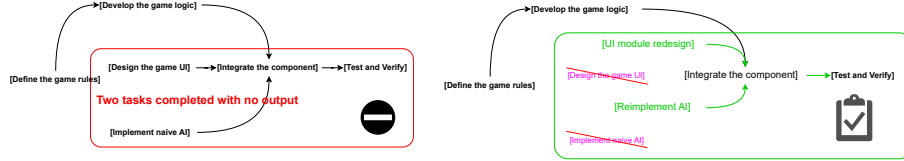
Fig. 3(b) illustrates a result of our dynamic updating process, where the system, upon receiving information about completed tasks, decides to add a bridging task to handle gaps and ensure the workflow continues smoothly.

Table 4: Success Rate (%) of Error handling with dynamically updating.

Task	Flow w/o Update	Flow
Website Design	46	<b>87</b>
Gobang Game	0	<b>93</b>
LaTeX Beamer Writing	67	<b>93</b>



(a) Conference website: no newly added subtask, only the workflow is updated.



(b) Gobang Game: bad subtasks exist, add two new subtasks for successfully completing this task.

Figure 3: Workflow and dynamic update in two cases.

We observed a significant difference in success rate between using dynamic updating and not, particularly in the Interactive Game section as shown in Table 4. The main issue arises when the previous agent fails to provide the necessary information, yet the second agent continues with its task, leading to a major disconnect in the code. This often results in Python being unable to compile due to missing or mismatched components. Similarly, in website design, the lack of required elements caused by this failure impacts the overall functionality and structure.

During the execution of subtasks, errors may arise due to the limitations of the LLM-based agent or under performance in certain tasks. Therefore, the ability of a multi-agent system to address such issues is essential.

**Limitations** Although we have generated multiple candidate workflows and selected the one with the highest modularity, it is still not the most efficient. With sufficient computing and data resources, a model trained specifically for workflow management could significantly enhance the system’s performance. For instance, the LLM could be designed to maximize a reward function centered on key performance indicators such as task completion speed, resource utilization, and minimizing disruptions in the workflow. Such training could lead to the development of more effective workflows. The workflow updater requires global information to function effectively, which can become problematic as the context length increases. This limitation could be addressed by employing a rig or a hierarchical approach to more precisely identify errors or areas lacking efficiency, thereby facilitating more targeted updates and improvements within the workflow.

## 6 CONCLUSION

We present Flow, a novel LLM-empowered multi-agent system that can dynamically adapt to unforeseen challenges for general tasks executions. With dynamically update the workflow by AOV graphs, our system has largely fulfilling the modularity requirements for completing complex tasks. We demonstrate our method through case studies on a series of experiments, ranging from website design, game development and LaTeX beamer creation as well as testing its capability on solving general benchmark tasks. Through objective evaluation metric and human feedback, we found Flow is to able to continuously enhance the flexibility during agent collaboration and thus significantly improve the execution efficiency with improved error tolerance and better performance.

## REFERENCES

- Carliss Y. Baldwin and Kim B. Clark. *Design Rules: The Power of Modularity Volume 1*. MIT Press, Cambridge, MA, USA, 1999. ISBN 0262024667.
- R.M. Belbin. *Team Roles at Work*. Butterworth-Heinemann, 2010. ISBN 9781856178006. URL <https://books.google.com.au/books?id=hF2yJzYfUBAC>.
- A. Bondy and U.S.R. Murty. *Graph Theory*. Graduate Texts in Mathematics. Springer London, 2011. ISBN 9781846289699. URL <https://books.google.com.au/books?id=HuDFMwZOwCsC>.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=EHg5GDnyq1>.
- T. DeMarco and T.R. Lister. *Peopleware: Productive Projects and Teams*. Addison-Wesley, 2013. ISBN 9780321934116. URL <https://books.google.com.au/books?id=DV1sAQAAQBAJ>.
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. MetaGPT: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=VtmBAGCN7o>.
- Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for "mind" exploration of large language model society. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Zhiwei Liu, Weiran Yao, Jianguo Zhang, Le Xue, Shelby Heinecke, Rithesh Murthy, Yihao Feng, Zeyuan Chen, Juan Carlos Nieves, Devansh Arpit, Ran Xu, Phil Mui, Huan Wang, Caiming Xiong, and Silvio Savarese. Bolaa: Benchmarking and orchestrating llm-augmented autonomous agents, 2023. URL <https://arxiv.org/abs/2308.05960>.
- J.J. Moder, C.R. Phillips, and E.W. Davis. *Project Management with CPM, PERT, and Precedence Diagramming*. Van Nostrand Reinhold, 1983. ISBN 9780442254155. URL <https://books.google.com.au/books?id=WmhRAAAAMAAJ>.
- OpenAI. Gpt-4o mini: Advancing cost-efficient intelligence. <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>, 2024. Accessed: 2024-09-29.
- Archiki Prasad, Alexander Koller, Mareike Hartmann, Peter Clark, Ashish Sabharwal, Mohit Bansal, and Tushar Khot. Adapt: As-needed decomposition and planning with language models. *arXiv*, 2023.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation, 2021. URL <https://arxiv.org/abs/2102.12092>.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. Reflexion: language agents with verbal reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=vAE1hFckW6>.
- Significant Gravitas. AutoGPT. <https://github.com/Significant-Gravitas/AutoGPT>. MIT License.

- Hongda Sun, Weikai Xu, Wei Liu, Jian Luan, Bin Wang, Shuo Shang, Ji-Rong Wen, and Rui Yan. DetermLR: Augmenting LLM-based logical reasoning from indeterminacy to determinacy. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9828–9862, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.531. URL <https://aclanthology.org/2024.acl-long.531>.
- H.A. Taha. *Operations Research an Introduction*. Pearson, 2017. ISBN 9780134444017. URL <https://books.google.com.au/books?id=HbpKjwEACAAJ>.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models, 2023. URL <https://arxiv.org/abs/2305.16291>.
- Yaoxiang Wang, Zhiyong Wu, Junfeng Yao, and Jinsong Su. Tdag: A multi-agent framework based on dynamic task decomposition and agent generation, 2024. URL <https://arxiv.org/abs/2402.10178>.
- Michael Wooldridge and Nicholas R. Jennings. Pitfalls of agent-oriented development. In *Proceedings of the Second International Conference on Autonomous Agents*, AGENTS ’98, pp. 385–391, New York, NY, USA, 1998. Association for Computing Machinery. ISBN 0897919831. doi: 10.1145/280765.280867. URL <https://doi.org/10.1145/280765.280867>.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversation. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2024.
- 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*, 2023. URL [https://openreview.net/forum?id=WE\\_vluYUL-X](https://openreview.net/forum?id=WE_vluYUL-X).
- Yining Ye, Xin Cong, Shizuo Tian, Yujia Qin, Chong Liu, Yankai Lin, Zhiyuan Liu, and Maosong Sun. Rational decision-making agent with internalized utility judgment, 2024. URL <https://openreview.net/forum?id=1lpNNQSzZv>.
- Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jialong Wu, Tiannan Wang, Shi Qiu, Jintian Zhang, Jing Chen, Ruipu Wu, Shuai Wang, Shiding Zhu, Jiyu Chen, Wentao Zhang, Xiangru Tang, Ningyu Zhang, Huajun Chen, Peng Cui, and Mrinmaya Sachan. Agents: An open-source framework for autonomous language agents. 2023. URL <https://arxiv.org/abs/2309.07870>.

# APPENDIX

## CONTENTS

<b>A Human Evaluation Process</b>	<b>14</b>
<b>B Experiment setups</b>	<b>14</b>
B.1 Experiment setup: LaTeX Generation . . . . .	14
B.2 Experiment setup: Go-bang game . . . . .	16
B.3 Experiment setup: Website Generation . . . . .	17
<b>C Examples of Flow’s Workflow</b>	<b>17</b>
C.1 Example workflow . . . . .	19
C.2 Pseudocode for updating AOV . . . . .	20

## A HUMAN EVALUATION PROCESS

Sometimes, LLMs can correctly fulfill each requirement of a task, but the quality of completion may vary. In such cases, human evaluation is necessary to assess the quality of the output. For each task, the final output of each Multi-Agent framework was evaluated by 50 participants, who ranked the outputs from best to worst. Points were awarded based on the rankings, with 1st place receiving 4 points and 2nd place receiving 3 points and so on. The final result was determined by calculating the average score. The detail distribution is shown in Fig. 5

## B EXPERIMENT SETUPS

### B.1 EXPERIMENT SETUP: LATEX GENERATION

**User input**

I am a lecturer teaching a machine learning course to research students, I am preparing lecture slides on various reinforcement learning algorithms.

Note that:

- 1). Given that the lecture duration is 2 hours, the slides should span approximately 30 pages.
- 2). For each reinforcement learning algorithm covered, the slides will include the following key components: the motivation behind the algorithm, the problem it aims to solve, an intuitive solution, and the detailed mathematical equations that underpin the method.
- 3). It is essential that the lecture is comprehensive and self-contained, providing students with a clear understanding of each algorithm from both a conceptual and technical perspective.

The task involves generating a LaTeX Beamer presentation, which is a popular LaTeX class used for creating professional-quality slides with various templates and effects. In this experiment, the objective is to produce presentations with different configurations, assessing the system’s ability to follow instructions. The experiment includes the following configurations:

- Config 1: A 30-slide presentation, including motivation, problem statement, intuitive solution, and detailed mathematical equations.
- Config 2: A 20-slide presentation, including motivation, problem statement, intuitive solution, and detailed mathematical equations.
- Config 3: A 30-slide presentation, including motivation, problem statement, intuitive solution, and pseudocode.
- Config 4: A 20-slide presentation, including only motivation and intuitive solution.
- Config 5: A 30-slide presentation, including motivation, problem statement, intuitive solution, and detailed mathematical equations.

The goal is to examine the system’s ability to follow specific instructions while generating over 20 and 30 slides in different scenarios.

This task is well-suited for evaluation because it requires not only text generation but also an understanding of formatting and presentation logic. It serves as a comprehensive test of multitasking and reasoning capabilities. The structured nature of LaTeX allows for a rigorous assessment of the agent’s ability to manage complex, multi-component tasks, thereby highlighting the strengths of our method.

**Evaluation Metrics:** The following metrics are used to assess the performance of the generated LaTeX Beamer presentations:

- (1) **Compilable:** Verifies whether the LaTeX code compiles into a valid Beamer presentation. A successful compilation is rewarded with a score of 1, otherwise 0.
- (2) **Completeness:** Ensures that the final Beamer presentation includes all required components: motivation, problem, intuitive solution, and equations. Missing any of these results in a score of 0.



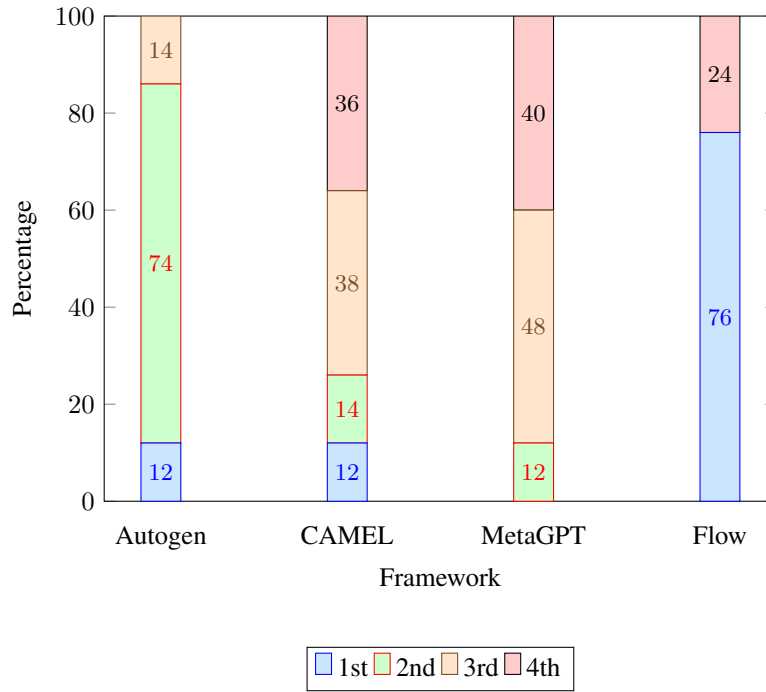


Figure 4: Ranking Distribution for conference website design. Shows that our results are better in the task of conference website design.

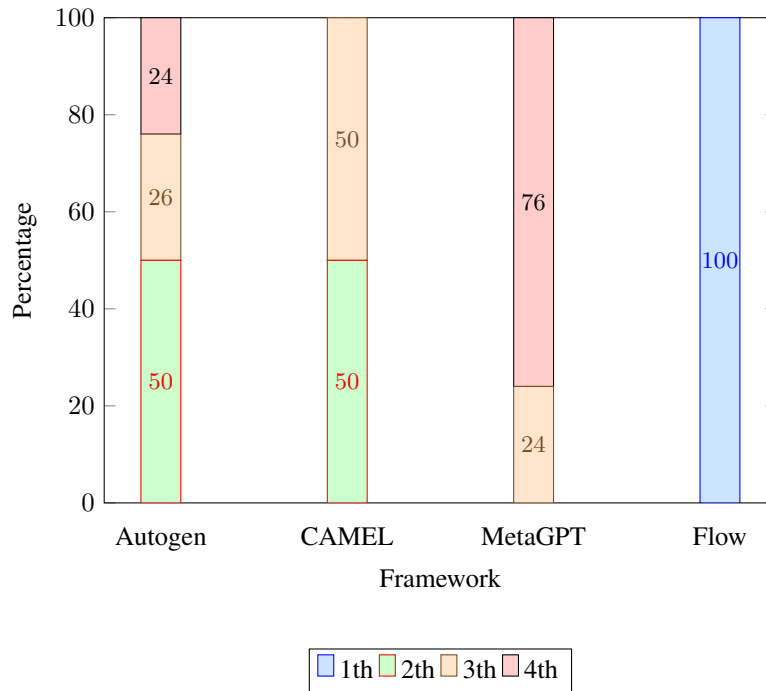


Figure 5: Ranking Distribution for Gobang game making. Shows that our results are better in the task of Gobang game making.

- (3) **Page Limit:** Assesses whether the presentation adheres to the specified page limits as outlined in the prompt.

The final result is calculated as the average of these three scores and shown in percentage.

Model	Compilable	Completeness	Page Limit	Overall Score
AutoGen Wu et al. (2024)	80	80	40	70
MetaGPT Hong et al. (2024)	80	80	20	60
CAMEL Li et al. (2023)	100	100	0	66
<b>Ours</b>	100	100	80	93

Table 5: Average Scores of Each LLM based Multi-Agent framework on Latex Generation Over 5 Trials

## B.2 EXPERIMENT SETUP: GO-BANG GAME

User input
I am developing a Gobang game that includes a naive AI and a user interface. The game should end when either a player wins or the board is completely filled. The user interface must clearly indicate whose turn it is and display a message when the game concludes, specifying the winner. Additionally, the user should have the option to play as either black or white stones.

Gobang, also called Five in a Row, is a strategy board game where two players take turns placing black and white pieces on a grid. The objective is to be the first to align five consecutive pieces in a horizontal, vertical, or diagonal line. This experiment assesses our system’s ability to efficiently develop the game by utilizing parallelism to divide the development process into smaller, manageable tasks, such as game logic, AI move generation, and user interface (UI) design. We apply the same approach, taking the average score from five trials.

**Evaluation Metrics:** The following metrics are used to assess the performance of the generated Gobang game:

- (1) **Compilable:** The code compiles without errors. Any error that causes a termination will result in a score of 0.
- (2) **Interactable:** Properly supports both user and AI moves. If both functions are achieved score 1 else 0
- (3) **Game Rule:** Ends correctly when five pieces are aligned, correct terminated will result in 1 final score.

Each of these metrics is scored as 0 or 1, and the final result is calculated as the average of these scores and turn into percentage. These metrics allow for a comprehensive assessment of the efficiency, accuracy, and adaptability of each framework in developing a functional Gobang game with AI capabilities.

Model	Compilable	Interactable	Game Rule	Overall Score
AutoGen Wu et al. (2024)	80	60	40	60
MetaGPT Hong et al. (2024)	100	100	20	73
CAMEL Li et al. (2023)	40	40	0	27
<b>Ours</b>	100	100	100	100

Table 6: Average Scores of Each LLM based Multi-Agent framework on Gobang Game Over 5 Trials

### B.3 EXPERIMENT SETUP: WEBSITE GENERATION

User input
<p>I am designing a website for the International Conference on Learning Representations (ICLR2025), which will take place from April 27, 2025, to May 1, 2025, in San Francisco, California, United States. The conference is organized by the International Association for Learning Representations.</p> <p>Note that:</p> <ol style="list-style-type: none"> <li>1). For each section, I would like to see example HTML content. Additionally, a sample CSS stylesheet should be provided to style the website. The content must be professional, clear, and appropriate for an international academic conference.</li> <li>2). The website should include all the provided details, including a comprehensive conference schedule and a section dedicated to the conference venue, featuring a map.</li> </ol>

We tasked the systems with developing a comprehensive website for the ICLR conference to evaluate their ability to handle complex tasks that require both flexible task coordination and effective problem-solving. This task tested the systems' ability to manage multiple interdependent steps, such as designing user interfaces, ensuring functionality, and adhering to specific design guidelines.

**Evaluation Metrics:** The following metrics are used to assess the performance of the generated website:

- (1) **Compilable:** Checks if the HTML renders into a functioning website. If yes then score 1, can't render will result of score 0
- (2) **Basic Information:** Verifies the presence of essential details like conference name, date, location, and organizer. Missing any of the information will caused the score to be 0
- (3) **Sections:** Ensures inclusion of all required sections, with a focus on the schedule and venue as prompt asked. Missing required part in prompt will result of 0 in score.

By presenting a real-world scenario involving intricate requirements, we were able to observe how well the systems could break down a large project into manageable components and coordinate efforts across different tasks.

Model	Compilable	Basic Information	Sections	Overall Score
AutoGen <a href="#">Wu et al. (2024)</a>	80	80	60	73
MetaGPT <a href="#">Hong et al. (2024)</a>	100	100	40	80
CAMEL <a href="#">Li et al. (2023)</a>	80	80	0	53
<b>Ours</b>	80	80	80	80

Table 7: Average Scores of Each LLM based Multi-Agent framework on Website Design Over 5 Trials

## C EXAMPLES OF FLOW'S WORKFLOW

In this section, we present examples of the actual workflows generated by Flow.

Fig.6 showing Flow's workflow in generating LaTeX Beamer, Flowconcurrently generates the four required components for each algorithm: motivation, problem, intuitive solution, and mathematical equations.

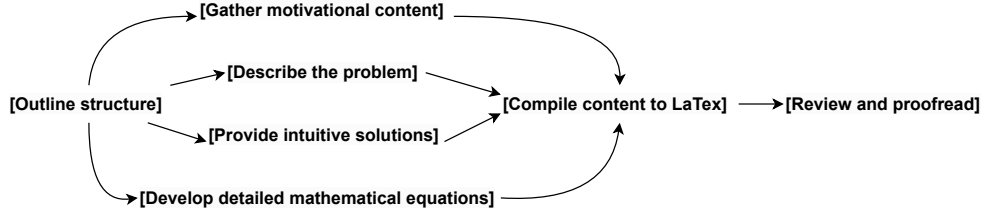


Figure 6: Workflow of LaTeX Beamer Writing in Flow

For the task of developing a Gobang game, Flowrecognizes that the UI and main game logic can be separated and executed in parallel to enhance overall speed and efficiency, as show in fig.7. Additionally, there remains a clear sequential process; for instance, the game rules must be defined first before the corresponding code can be deployed.

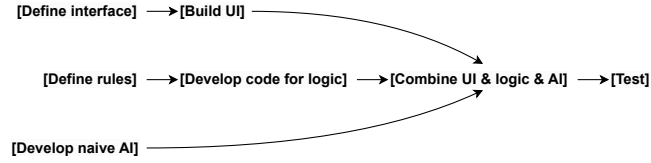


Figure 7: Workflow of Gobang game generation

For the task of generating a website show in Fig.8, Flowtreats different parts of the HTML as individual subtasks, which helps to increase overall speed. Additionally, dividing the process into separate components allows for parallel execution and improved modularity, ensuring that if an issue arises in one part of the HTML, it will not impact the performance of other sections. This approach enhances both efficiency and fault tolerance.

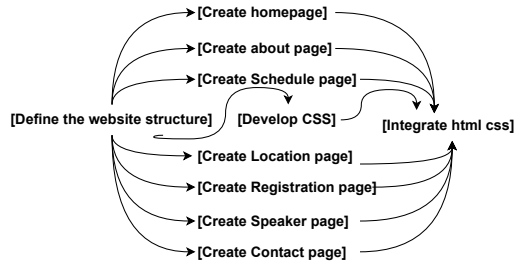


Figure 8: Workflow of Conference Website Generation

## C.1 EXAMPLE WORKFLOW

```

"task0": {
  "params": {
    "task": "Define the structure and layout for the LaTeX document, including title page, table of contents, and section headers",
    "agent": 0
  },
  "next": [
    "task3"
  ],
  "prev": [],
  "remaining_dependencies": 0,
  "status": "completed",
  "data": "Here is a proposed structure and layout for a LaTeX document that encompasses the requirements for a two-hour lecture on different reinforcement l
},
"task1": {
  "params": {
    "task": "Gather and compile information on various reinforcement learning algorithms, including their motivations, problems, and intuitive solutions",
    "agent": 1
  },
  "next": [
    "task3",
    "task4",
    "task5"
  ],
  "prev": [],
  "remaining_dependencies": 0,
  "status": "completed",
  "data": "### Compiled Information on Reinforcement Learning Algorithms\n\n### 1. Q-Learning\n\n**Motivation**: To develop a model-free reinforcement learn
),

```

Figure 9: A workflow of website design in VSCode

Q-Learning	Problem Statement: Q-Learning	Intuitive Solution: Q-Learning	Mathematical Foundations: Q-Learning
<b>Motivation</b> <ul style="list-style-type: none"> <li>Discover the optimal action-value function.</li> <li>Learn optimal policies in an off-policy manner.</li> </ul> <b>Key Insight</b> <ul style="list-style-type: none"> <li>Learn from the consequences of actions rather than the specific environment model.</li> </ul>	<ul style="list-style-type: none"> <li>How can we construct an agent that learns the optimal action-value function from arbitrary exploration of the environment?</li> <li>Challenge: Ensuring convergence to the optimal policy in stochastic environments.</li> </ul>	<ul style="list-style-type: none"> <li>Update Q-values based on observed rewards and estimated future rewards.</li> <li>Use the Bellman equation to iteratively improve Q-value estimates.</li> </ul>	<p>The Q-learning update rule is defined as:</p> $Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$ <p>This equation estimates the optimal action-value function despite the stochastic nature of the environment.</p>

(a): Flow

Q-Learning	Q-Learning: Problem	Q-Learning: Intuitive Solution	Q-Learning: Detailed Math
<b>Motivation:</b> <ul style="list-style-type: none"> <li>Understanding how agents learn optimal actions through trial and error.</li> </ul>	<ul style="list-style-type: none"> <li>Define the environment and states.</li> <li>Challenges in learning optimal policies.</li> </ul>	<ul style="list-style-type: none"> <li>Use of Q-values to represent action-value pairs.</li> <li>Update rule:</li> </ul> $Q(s, a) \leftarrow Q(s, a) + \alpha \left( r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right)$	<ul style="list-style-type: none"> <li>Explanation of terms: <ul style="list-style-type: none"> <li><math>s</math>: current state</li> <li><math>a</math>: action taken</li> <li><math>r</math>: reward received</li> <li><math>\alpha</math>: learning rate</li> <li><math>\gamma</math>: discount factor</li> </ul> </li> </ul>

(b): AutoGen

Motivation	Mathematical Details	Algorithm Steps	algorithms/algorithm1.tex
<b>The motivation behind Algorithm 1</b> is to provide a robust solution for reinforcement learning tasks by optimizing the decision-making process through exploration and exploitation strategies. <ul style="list-style-type: none"> <li>Efficient exploration of the state space.</li> <li>Balancing exploration and exploitation.</li> <li>Adaptability to dynamic environments.</li> </ul>	<p>The core of Algorithm 1 is based on the following mathematical formulation:</p> $Q(s, a) \leftarrow Q(s, a) + \alpha \left( r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right) \quad (1)$ <p>Where:</p> <ul style="list-style-type: none"> <li><math>Q(s, a)</math> is the action-value function.</li> <li><math>\alpha</math> is the learning rate.</li> <li><math>r</math> is the reward received after taking action <math>a</math> in state <math>s</math>.</li> <li><math>\gamma</math> is the discount factor.</li> <li><math>s'</math> is the next state after taking action <math>a</math>.</li> </ul>	<p>The steps of Algorithm 1 are as follows:</p> <ol style="list-style-type: none"> <li>Initialize the action-value function <math>Q(s, a)</math> arbitrarily.</li> <li>For each episode: <ul style="list-style-type: none"> <li>Initialize the starting state <math>s</math>.</li> <li>For each step of the episode: <ul style="list-style-type: none"> <li>Choose an action <math>a</math> based on the current policy.</li> <li>Take action <math>a</math>, observe reward <math>r</math> and next state <math>s'</math>.</li> <li>Update the action-value function using the formula above.</li> <li>Set <math>s = s'</math>.</li> </ul> </li> </ul> </li> </ol>	

(c): Meta

Q-Learning Algorithm
<b>Introduction:</b> <ul style="list-style-type: none"> <li>Q-Learning is a model-free reinforcement learning algorithm used to learn the value of actions in a given state.</li> <li>It enables an agent to learn optimal policies without needing a model of the environment.</li> </ul> <b>Significance:</b> <ul style="list-style-type: none"> <li>Q-Learning is widely used due to its simplicity and effectiveness.</li> <li>It forms the foundation for many advanced reinforcement learning algorithms, including Deep Q-Networks (DQN).</li> <li>It allows agents to learn from past experiences and improve their decision-making over time.</li> </ul> <b>Q Value Function:</b> <ul style="list-style-type: none"> <li>The Q Value function, denoted as <math>Q(s, a)</math>, represents the expected cumulative reward of taking action <math>a</math> in state <math>s</math> and following the optimal policy thereafter.</li> </ul>

(d): CAMEL

Figure 10: Different Multi-Agent frameworks' LaTeX Beamer

## C.2 PSEUDOCODE FOR UPDATING AOV

**Algorithm 1:** Flow**Data:** Task Requirements  $T$ , Initialization Prompt  $P_{\text{init}}$ , Update Prompt  $P_{\text{update}}$ **Result:** Optimized Multi-Agent Workflow

---

```

// Step 1: Implement a Workflow using a dictionary structure
1 Initialize workflow formulation by defining the task dictionary  $\tilde{G}$  where each key  $v \in V$  maps
  to a dictionary containing:  $\tilde{G}[v] = \{\text{status}, \text{data}, \text{num\_parents\_not\_completed}, \text{child}, \text{agent}\}$ 
// Step 2: Generate an Initial Workflow
2  $\tilde{G} \leftarrow \text{UpdateGraph}(\{\}, P_{\text{init}}, T)$ ;
// Step 3: Workflow Refinement and Dynamic Updating
3 while there exists at least one sub-task in  $\tilde{G}$  that is not completed do
4   if an update to the workflow is required then
5     // Generate and Select the Best Updated Workflow
6      $\tilde{G} \leftarrow \text{UpdateGraph}(\tilde{G}, P_{\text{update}}, T)$ ;
7     Update workflow dictionary  $\tilde{G}$  to  $\tilde{G}_{\text{best}}$ ;
8     // Regenerate Execution Plan and Reallocate Agents
9     Perform Topological Sort on  $\tilde{G}$  to obtain updated execution order  $\sigma$ ;
10    Assign agents  $A_j$  to their respective sub-tasks  $\mathcal{T}_j \subseteq V$ ;
11  end
12  // Execute Available Sub-tasks in Parallel
13  foreach sub-task  $v_i \in V$  do
14    if status of  $v_i$  is not started and  $\tilde{G}[v_i].\text{num\_parents\_not\_completed} == 0$  then
15      if agent  $a_j$  is available then
16        | Assign agent  $a_j$  to sub-task  $v_i$ ;
17      else
18        | Clone agent  $a'_j$ ;
19        | Assign cloned agent  $a'_j$  to sub-task  $v_i$ ;
20      end
21      // Execute sub-task  $v_i$  in parallel
22      Execute  $v_i$  using agent  $a_j$  or cloned agent  $a'_j$  concurrently;
23      // Update Sub-task Status and Data
24      Update status of sub-task  $v_i$  to in progress;
25      // After execution, update related data
26      Update output of sub-task  $v_i$  to  $\tilde{G}[v_i].\text{data}$ ;
27       $\tilde{G}[v_i].\text{status} \leftarrow \text{"completed"}$ ;
28      // Update Child Tasks' Parent Completion Count
29      foreach child task  $c \in \tilde{G}[v_i].\text{child}$  do
30        |  $\tilde{G}[c].\text{num\_parents\_not\_completed} \leftarrow \tilde{G}[c].\text{num\_parents\_not\_completed} - 1$ ;
31      end
32    end
33  end

```

---



**Algorithm 2:** Helper Function for Updating Graph

---

```

1080
1081
1082 1 Function UpdateGraph( $\tilde{G}, P, T$ ):
1083   // Generate updated candidate workflows using LLM
1084 2    $\{\tilde{G}_1, \tilde{G}_2, \dots, \tilde{G}_K\} \leftarrow f(\tilde{G}, P, T)$ ;
1085   // Initialize selection variables
1086 3    $P_{\max} \leftarrow -\infty$ ;
1087 4    $C_{\min} \leftarrow +\infty$ ;
1088 5    $\tilde{G}_{\text{optimal}} \leftarrow \text{None}$ ;
1089   // Evaluate each candidate workflow
1090 6   for each candidate workflow  $\tilde{G}_k$  in  $\{\tilde{G}_1, \tilde{G}_2, \dots, \tilde{G}_K\}$  do
1091 7     Compute Parallelism  $P_k \leftarrow P_{\text{avg}}(\tilde{G}_k)$ ;
1092 8     Compute Dependency Complexity  $C_k \leftarrow C_{\text{dependency}}(\tilde{G}_k)$ ;
1093 9     if  $P_k > P_{\max}$  or ( $P_k == P_{\max}$  and  $C_k < C_{\min}$ ) then
1094 10        $P_{\max} \leftarrow P_k$ ;
1095 11        $C_{\min} \leftarrow C_k$ ;
1096 12        $\tilde{G}_{\text{optimal}} \leftarrow \tilde{G}_k$ ;
1097
1098 13 return  $\tilde{G}_{\text{optimal}}$ ;

```

---