

ADAPTIVE IN-CONVERSATION TEAM BUILDING FOR LANGUAGE MODEL AGENTS

Anonymous authors

Paper under double-blind review

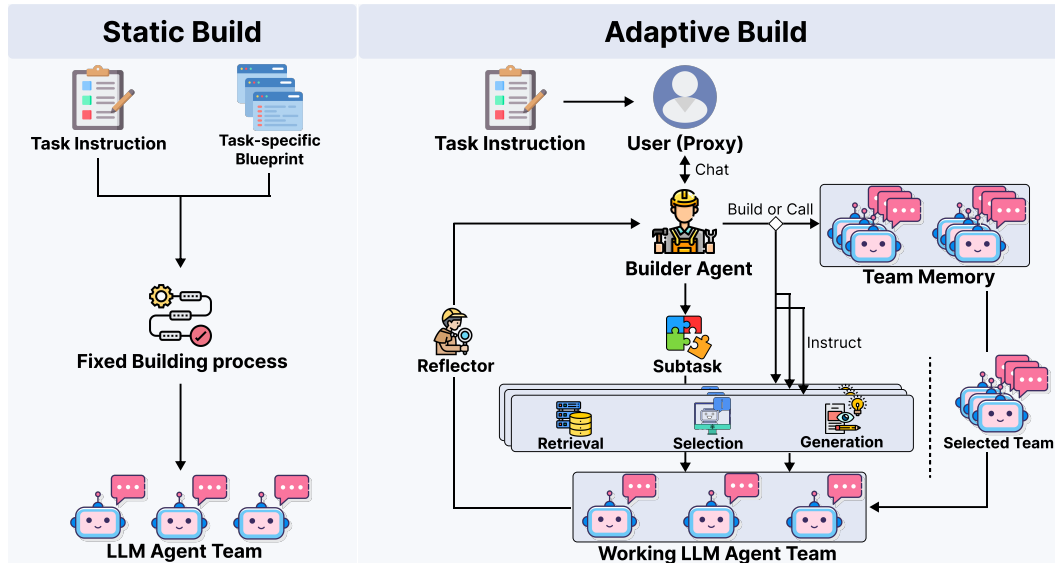


Figure 1: Two team-building paradigms for LLM agents. The "Static Build" paradigm (Chen et al., 2023; Suzgun & Kalai, 2024a; Chen et al., 2024) statically builds a team according to task instructions and a building blueprint through a fixed process. Our "Adaptive Build" paradigm uses an adaptive builder agent to form different teams during the task-solving procedure without extra instructions.

ABSTRACT

Leveraging multiple large language model (LLM) agents has shown to be a promising approach for tackling complex tasks, while the effective design of multiple agents for a particular application remains an art. It is thus intriguing to answer a critical question: *Given a task, how can we build a team of LLM agents to solve it effectively?* Our new adaptive team-building paradigm offers a flexible solution, realized through a novel agent design named *Captain Agent*. It dynamically forms and manages teams for each step of a task-solving process, utilizing nested group conversations and reflection to ensure diverse expertise and prevent stereotypical outputs, allowing for a flexible yet structured approach to problem-solving. A comprehensive evaluation across six real-world scenarios demonstrates that *Captain Agent* significantly outperforms existing multi-agent methods with 21.94% improvement in average accuracy, providing outstanding performance without requiring task-specific prompt engineering. Our exploration of different backbone LLM and cost analysis further shows that *Captain Agent* can improve the conversation quality of weak LLM and achieve competitive performance with extremely low cost, which illuminates the application of multi-agent systems.

1 INTRODUCTION

The success of large language model (LLM) agents (Yao et al., 2022; Yang et al., 2023a; Furuta et al., 2024; Yang et al., 2024a; Hong et al., 2024) with its outstanding in-context learning (Dong

054 et al., 2022; Brown et al., 2020; Yang et al., 2023b; Dai et al., 2023; Li et al., 2023b), planning (Sun
055 et al., 2024; Xie et al., 2024; Liu et al., 2023a; Valmeekam et al., 2022; Wei et al., 2022a; Yuan
056 et al., 2023b; Zheng et al., 2024), tool-using (Qin et al., 2023a;b; Schick et al., 2024; Cai et al.,
057 2023; Yuan et al., 2023a; Paranjape et al., 2023; Zhang et al., 2024b; Huang et al., 2023; Ma et al.,
058 2024), and conversation (Fernandes et al., 2023; Wang et al., 2023c; Yang et al., 2024b) capabilities
059 allow us to relate human’s team building and collaboration abilities to the multiple language model
060 agents (multi-agent) system (Wang et al., 2023a; Xi et al., 2023; Wu et al., 2023; Suzgun & Kalai,
061 2024a; Hong et al., 2023; Zhang et al., 2024b; 2023a; Valmeekam et al., 2023; Wang et al., 2024;
062 Saha et al., 2023; Liang et al., 2023; Du et al., 2023; Chen et al., 2024). Humans have developed
063 abilities that enable us to form teams and effectively solve problems. These abilities are rooted
064 in communication, social cognition, problem-solving and decision-making, social learning and
065 imitation, and shared intentionality (Elimari & Lafargue, 2020; Confer et al., 2010). The interplay of
066 the above abilities allows people to organize different teams for problems to ensure that tasks are
067 completed successfully, which brings us to a critical question in a multi-agent system:

068 *Given a task, how can we build a team of LLM agents to solve it effectively?*

069 A straightforward paradigm would be to build a static agent team beforehand based on the task
070 instruction and let them solve the task collaboratively (Chen et al., 2023; Wu et al., 2023). However,
071 this *static build* method necessitates maintaining a team with all the required expertise for the whole
072 task cycle. As the complexity of the task increases, the total number of team members may grow
073 significantly. Always proceeding with such a large team makes it challenging to manage the team
074 members effectively and efficiently. Furthermore, static teams may lack the adaptability to respond
075 to dynamic changes in task requirements or unforeseen challenges. Imagine a prehistoric human
076 tribe: was everyone involved in every task? The answer is unlikely affirmative. Those responsible
077 for hunting may not participate in medical care and those responsible for cooking may not involve
078 themselves in management. The major task, survival, was ensured by each individual group sticking
079 to their roles and subtasks. In fact, when human organizations handle a complex task, we tend to
080 form multiple teams for each subtask at different stages of the task-solving procedure, which still
081 guarantees a diverse set of expertise is leveraged demanded by the task complexity (Mao et al., 2016).

082 Inspired by how humans assemble teams for a complex task, we introduce a new multi-agent
083 team-building paradigm: *adaptive build*. This paradigm facilitates the flexible assembly of agents
084 with specific skills and knowledge as demands evolve in the process of task-solving. To realize
085 this paradigm, we propose a new adaptive builder agent, Captain Agent, to build, manage, and
086 maintain agent teams for each problem-solving step in the conversation. Captain Agent has two core
087 components: (1) adaptive multi-agent team building and (2) nested group conversation and reflection.
088 Captain Agent will communicate with a User Proxy, who can provide the general task instructions at
089 the beginning. When assigned a task, Captain Agent begins by formulating a strategic plan. This plan
090 involves a cyclical process that continues until the task is successfully completed. In the first phase of
091 the cycle, Captain Agent identifies a specific subtask, outlines the necessary roles, and assembles a
092 team of agents equipped with the appropriate tools. In the subsequent phase, this team engages in
093 a dialogue with a versatile tool to address the subtask. Upon completion, a reflector LLM reviews
094 the process and provides Captain Agent with a detailed reflection report. Based on this feedback,
095 Captain Agent either adjusts the team composition or the subtask instructions and repeats the cycle
or concludes the task and presents the final outcomes.

096 We evaluate state-of-the-art multi-agent approaches for complex task solving and our adaptive build
097 approach with Captain Agent on six real-world scenarios, including many mathematics problem-
098 solving (Hendrycks et al., 2021b), data analysis (Hu et al., 2024b), programming (Le et al., 2020),
099 scientific problem-solving (Wang et al., 2023b) (Physics and Chemistry), and world-information
100 retrieval (Mialon et al., 2024). Our experimental results demonstrated the outstanding ability of
101 Captain Agent in various scenarios without heavy prompt engineering for each scenario but only
102 the basic instructions. Captain Agent achieves distinguishing results compared to other single and
103 multi-agent methods and frameworks when using the same prompt for each task, with an average of
104 21.94% improvement on average accuracy. Ablation studies on static and adaptive building paradigms
105 show that the adaptive team outperforms the static team in four of five scenarios (and matches in one
106 scenario), exhibiting the superiority of the adaptive build paradigm across different scenarios. We
107 also demonstrated that handcraft agents and handcraft tools contribute equally to the final results. We
further explore the influence of different backbone LLM for both Captain Agent and nested group

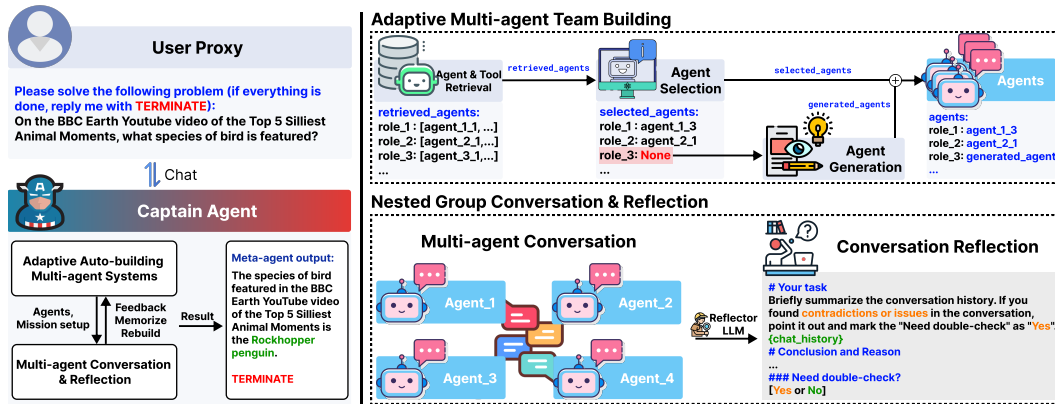


Figure 2: The overall workflow of Captain Agent is: given a user instruction, Captain Agent will plan the task, build an agent team from retrieval and generation, and let the agents solve a decomposed, planned task collaboratively in a group chat. A reflection LLM will review and report the conversation history to Captain Agent. Captain Agent will then conclude or continue solving the problem with a modified team and instructions.

chat members or only for nested group chat members. We observe that: (1) Captain Agent with a strong backbone can improve the quality of nested group chat in which the members equipped with weak backbone, and (2) a small model with distinguishable instruction following ability can achieve outstanding performance with low cost.

2 ADAPTIVE IN-CONVERSATION TEAM BUILDING

The proposed Captain Agent contains two key components: (1) adaptive multi-agent team-building, which involves agent and tool retrieval, selection, and generation, and (2) nested group conversation with a reflection mechanism within the multi-agent system.

2.1 OVERVIEW

The overall workflow of Captain Agent is illustrated in Figure 2. Given a task, Captain Agent is prompted to derive a plan before task execution. According to the plan, Captain Agent will repeat the following two steps until it thinks the task is done and output the results: **(Step 1)** Captain Agent will first identify a subtask instructed by our prompt, list several roles needed for this subtask, and then create a team of agents accordingly by retrieval, selection, and generation. Each of these will be equipped with predefined tools retrieved from the tool library (Section 2.2); **(Step 2)** this team of agents will attempt to solve the subtask via conversation with the free-form tool using. Once it’s done, a reflector LLM will provide Captain Agent with a reflection report for it to decide whether to adjust the team or subtask instruction or to terminate and output the results (Section 2.3).

2.2 ADAPTIVE MULTI-AGENT TEAM BUILDING

After identifying a subtask in Step 1 following a corresponding prompt, Captain Agent will list several roles for the subtask. These roles will then pass into a retrieval, selection, and generation process guided by Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Gao et al., 2023; Ram et al., 2023). Created agents will be equipped with a well-designed profile (system message¹) and high-quality tools. We illustrated the whole process in Figure 3.

Agent and tool retrieval. Captain Agent will prompt n required roles $\{r_i | i \in 1, \dots, n\}$ with detailed descriptions, including required skills and a possible role name. We use "expert" in Captain Agent prompt to make this process natural. We then retrieve top- k_1 agents and top- k_2 tools according to the sentence embedding similarity between the role’s description and the agent/tool description recorded

¹System message is used to define an agent’s persona and task-specific instructions.

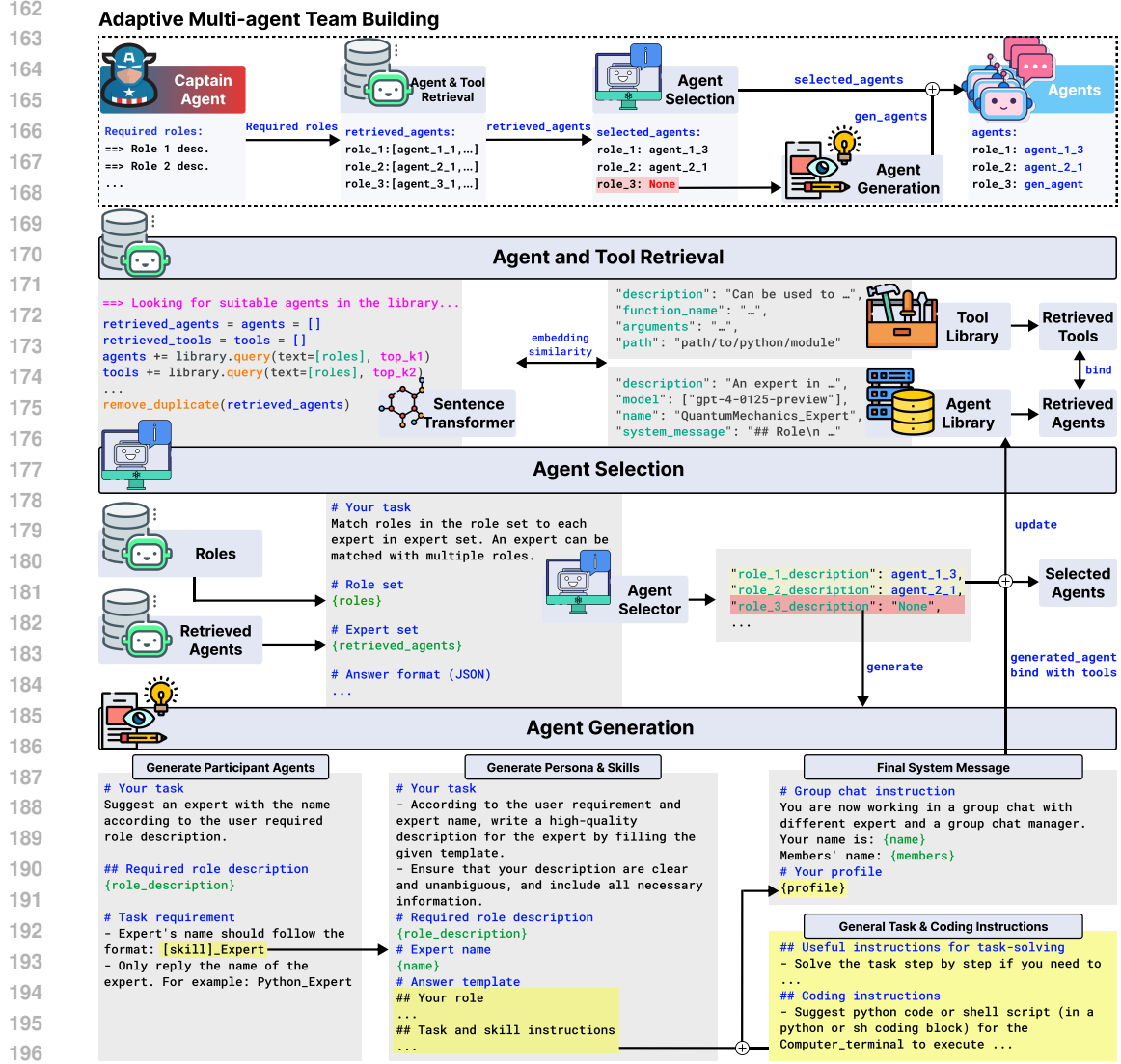


Figure 3: Workflow for adaptive multi-agent team building. We retrieve candidate agents and tools according to the roles’ description prompted by Captain Agent. Candidate agents and tools will further be linked to a role under the advice of the agent selector. If no agent is linked to a role, a generate process will be performed to create a new agent. It will generate the agent’s name and task-specific instructions, combined with general task and coding skills and group chat instructions as the final system message.

in the library. We use Sentence Transformer to calculate the embedding for description between the role and library agents/tools and use cosine similarity as the metric to evaluate the similarity between two sentences, as follows:

$$\text{top-}k_1 \text{ CosineSimilarity}(f(r_i), f(a_{\text{lib}})) \rightarrow \text{RetrievedAgents}, \quad (1)$$

$$\text{top-}k_2 \text{ CosineSimilarity}(f(r_i), f(t_{\text{lib}})) \rightarrow \text{RetrievedTools}, \quad (2)$$

where k_1 and k_2 are the numbers of retrieved agents and tools from agent library a_{lib} and tool library t_{lib} , respectively, for i -th role r_i . $f(\cdot) \in \mathbb{R}^m$ denotes the sentence embedding extracted from a Sentence Transformer. After retrieval, each role will be assigned with k_1 agent candidates and k_2 valuable tools. We bind agent candidates with the retrieved tools by injecting the tool-using instruction into the corresponding agent’s system message.

Agent selection. We prompt an LLM-based agent selector to select the most suitable agent according to the role’s description given by Captain Agent and the retrieved agents’ description. A JSON

216 template is designed and provided for the agent selector to ensure the format is correct. Specifically,
217 we designed an abstention mechanism for the agent selector, in which the agent selector can output
218 "None" if there is no suitable agent for a role from the top- k_1 retrieved candidate list. This can prevent
219 irrelevant or redundant agents from being forced to be selected for the current task. The roles marked
220 with "None" will further go into the generation process described below.

221 **Agent generation.** We design an agent generation process for those roles with no linked agents at
222 the previous step. Specifically, we generate the agent's name and required skills according to the
223 role description given by Captain Agent. These instructions will be combined with general task and
224 coding instructions and group chat instructions as the final system message. We manually design the
225 general task and coding instructions, motivated by Chain-of-thought (CoT) (Wei et al., 2022b) and
226 Reflexion (Shinn et al., 2024). The final system message will also be compressed to a single-sentence
227 description, which is consumed by the nested group conversation (introduced in the next subsection).
228 We then retrieve tools from the tool library according to the description and inject the tool-using
229 instruction into the generated system message.

230 **Team Memory.** Once the team has been built, Captain Agent will cache it into its local memory
231 with a team name and each agent's detail, including name, system message, and the assigned tools.
232 Captain Agent can call the cached team anytime during the conversation with the user proxy. Calling
233 the cached team will not incur any API calls and thus will not introduce extra costs.

234 235 236 2.3 NESTED GROUP CONVERSATION AND REFLECTION

237 Agents selected and created in the adaptive multi-agent team-building process will join a nested
238 group chat room. They will be prompted to collect information from the user's task and solve a
239 subtask from Captain Agent by nested conversation. We then prompt a reflector LLM to retrieve and
240 review the conversation history and fill in the conclusion, the reason for the conclusion, possible
241 contradictions, and issues, and flag if the result needs a double check in the pre-designed template.

242 **Nested group conversation.** We perform nested group conversations by leveraging the AutoGen (Wu
243 et al., 2023) framework with a newly designed tool-using paradigm. AutoGen will put all agents
244 in a chat room and select the speaker for each turn by a group chat manager LLM according to the
245 conversation history and each agent's identity. A short description will be generated from the agent's
246 profile for the group chat manager. Agents' code and tool calling will be executed and fed back to the
247 conversation immediately. We inject the tool's description, path-to-python-module, and response case
248 into the related agent's system message. The agent can then write free-form code by following the
249 tools' description and path, naturally incorporating the tools into larger programs. Programs written
250 by all agents will be executed by a user proxy agent with a shared code execution environment, and
251 the results will be fed back to the conversation in real time.

252 **Conversation reflection.** The agent's output during the conversation can be inconsistent, including
253 factual errors, hallucinations, and stereotypes. Although other agents have a chance to adjust and
254 rectify this in conversation, they can also get stuck and cause problem-solving failure. Therefore,
255 we propose to detect such in-conversation contradictions and issues by prompting a reflector LLM
256 with a well-designed conversation summarizing prompt template. The reflector will flag the "need
257 double-check" as "Yes" when it detects such inconsistent content and provides a detailed reason. This
258 will trigger Captain Agent to start a verification process by constructing a new nested conversation to
259 double-check the previous results after receiving "Yes" on "need double-check."

260 261 2.4 BENEFITS OVER STATIC BUILD

262
263 A static team with a small number of team members may limit the team's ability coverage. Although
264 building a large number of agents with comprehensive persona or skill sets can address the limitation
265 in ability coverage, it is challenging for LLMs to handle a long context that introduces all the
266 participant members. Unexpectedly long contexts will primarily reduce the quality of the conversation.
267 Meanwhile, agents with redundant functionality will also be involved in the task-solving process. In
268 contrast, Captain Agent can adaptively select and build more optimized agent teams for the current
269 task, reducing the prompting load for LLMs and redundant output from irrelevant agents without
sacrificing the diversity in the agent team.

3 EVALUATION

3.1 EXPERIMENTAL SETUP

Table 1: Scenarios and the corresponding datasets we choose to perform our main experiments. We perform the main comparison experiments on the whole dataset except MATH. For MATH, we sampled a small subset according to the type distribution.

Scenario	Dataset	Size	Sample
Mathematics problems	MATH (Hendrycks et al., 2021a)	196	If $\frac{3x^2-4x+1}{x-1} = m$, and x can be any real number except 1, what real values can m NOT have?
Programming	HumanEval (Chen et al., 2021)	164	def truncate_number(number: float) ->float: """ Given a positive floating point number, it can be decomposed into and integer part (largest integer smaller than given number) and decimals (leftover part always smaller than 1). [Omitted] """
Data Analysis	DABench (Hu et al., 2024a)	257	Generate a new feature called "FamilySize" by summing the "SibSp" and "Parch" columns. Then, calculate the Pearson correlation coefficient (r) between the "FamilySize" and "Fare" columns.
World Information Retrieval	GAIA (Mialon et al., 2023)	165	On the BBC Earth YouTube video of the Top 5 Silliest Animal Moments, what species of bird is featured?
(Scientific) Chemistry	SciBench (Wang et al., 2023b)	41	Calculate the pressure in kilopascals exerted by 1.25 g of nitrogen gas in a flask of volume 250 cm ³ at 20°C.
(Scientific) Physics	SciBench (Wang et al., 2023b)	34	If the coefficient of static friction between the block and plane in the previous example is $\mu_s = 0.4$, at what angle θ will the block start sliding if it is initially at rest?

Scenarios and datasets. For evaluation, we select various real-world scenarios, including mathematics problem-solving, programming, data analysis, world information retrieval, and science problem-solving. Each scenario was chosen for its unique ability to demonstrate specific capabilities and performance metrics of the agent systems. This ensures a holistic assessment of Captain Agent against the baselines across various critical dimensions of computational and cognitive skills. We bind each scenario with a challenging open-source dataset, as shown in Table 1. Due to cost limitations, we sample a subset of MATH according to its original distribution of each question type.

Compared methods and implementation. For mathematics problems, programming, data analysis, and scientific scenarios, we investigate the performance of Captain Agent and four different methods, including Vanilla LLM (prompt an LLM once for an answer), AutoAgents (Chen et al., 2023), Meta-prompting (Suzgun & Kalai, 2024a), AgentVerse (Chen et al., 2024), DyLAN (Liu et al., 2023b), and a two-agent system (a system involving an Assistant agent with an Executor agent) realized with AutoGen (Wu et al., 2023). Specifically, we implement AutoAgents with AutoGen as the official implementation is unstable and unsuitable for large-scale experiments. For meta-prompting, we improve the code execution ability of meta-prompting by reproducing it with the AutoGen framework. All these methods are equipped with a gpt-4-0125-preview backbone and use the same task-specific prompt (refer to Appendix E).

For world information retrieval scenarios, we compare Captain Agent with the top-5 baselines (with reference) reported to the GAIA validation leaderboard, which includes AutoGen: GAIA_Orchestrator (a specific three-agent setting organized by an Orchestrator agent designed for GAIA) (GAIA_Orchestrator, 2024), FRIDAY (Wu et al., 2024), Warm-up Act², and HuggingFace Agent (Huggingface, 2024). All these baselines have a gpt-4-1106-preview backbone, except the HuggingFace Agent equipped with an LLaMA-3-70B as the backbone.

For Captain Agent, we adopt all-mpnet-base-v2 to calculate the sentence embedding for agent and tool retrieval. A User Proxy Agent will communicate with Captain Agent by providing the feedback of code execution, tool calling (adaptive build), nested conversation reflection results, and a default reply: *I'm a proxy, and I can only execute your code and tool or end the conversation. If you think the problem is solved, please reply to me only with 'TERMINATE.'*

Agent and tool library. We initialize our agent library based on a small subset of problem instances from each dataset (~20 questions per dataset described in Section 3.4) in Table 1. Specifically, we run Captain Agent on the subset and iteratively update the library by adding the generated agents and keeping our agent library unchanged during the main experiment. Our agent library also supports all hand-crafted agents (of the ConversableAgent class) archived in AutoGen (details in Appendix G). All these agents follow the ConversableAgent interface to converse with each other.

²Warm-up Act has no official implementation.

Table 2: Comparison results on different real-world scenarios. We record each scenario’s accuracy for each baseline and Captain Agent, and mark the best results in **bold**. We adopt gpt-4-0125-preview as the backbone LLM model for all baselines and Captain Agent.

Method	Mathematics	Programming	Data Analysis	(Sci) Chemistry	(Sci) Physics	Avg.
Vanilla LLM	51.53	84.76	6.61	39.02	31.25	40.98
Meta-prompting	68.88	19.51	39.69	41.46	43.75	43.47
AutoAgents	56.12	84.76	57.98	60.98	50.00	63.58
DyLAN	62.24	90.24	-	45.45	51.16	-
AgentVerse	69.38	42.68	-	42.42	37.21	-
AutoGen: Assistant + Executor	74.49	93.90	82.88	60.98	43.75	79.89
Captain Agent	77.55	96.95	88.32	65.85	53.12	84.25

Table 3: Comparison results on world-information retrieval scenario (GAIA validation). We report the accuracy at each level and the average accuracy over three levels and mark the best results in **bold**. Captain Agent achieves the best with minimal prompt engineering.

Method	Level 1	Level 2	Level 3	Avg.
AutoGPT4	13.21	0.00	3.85	4.85
GPT4 Turbo	20.75	5.81	0.00	9.70
GPT4 + manually selected plugins	30.30	9.70	0.00	14.6
Captain Agent (Llama-3-70B-Instruct)	28.30	11.63	0.00	15.15
Huggingface-Agent (Llama-3-70B-Instruct)	30.19	11.63	7.69	16.97
Warm-up Act	35.19	15.12	0	17.58
Captain Agent (gpt-4o-mini)	32.08	16.27	3.85	19.39
FRIDAY	45.28	34.88	11.54	34.55
AutoGen: GAIA_Orchestrator	54.72	38.31	11.54	39.39
Captain Agent (gpt-4-0125-preview)	56.60	39.53	11.54	40.60

Our tool library consists of a suite of callable Python functions intended for freeform coding. The agents can freely import functions from the tool library and write free-form code to integrate the outputs to handle sophisticated tasks (see also Appendix F and H). The library contains three main categories of tools: math, data analysis, and world information retrieval. For each category, we summarize the patterns of the corresponding dataset and manually craft a set of functions that suit the tasks.

3.2 EVALUATION PROTOCOL

For mathematics, data analysis, and science scenarios, we report the accuracy of each method by comparing the final result from each method and ground truth. To ensure fairness in evaluation, we transform different result formats into a uniform format, preventing the correct answer from being judged incorrect due to format mismatches. For programming scenarios, we run the code provided from each method and output a unique token if the code successfully passes all tests. We then count the success token and calculate the accuracy for each method.

3.3 MAIN RESULTS

Table 2 and 3 report the comparison results between Captain Agent and eight different baselines on six real-world scenarios. Baseline results on world information retrieval are extracted directly from the GAIA leaderboard.

Findings 1: Diverse agents can help trigger accurate expertise output for problem-solving.

By comparing the results from Captain Agent, AutoAgents, and AutoGen Assistant + Executor, we observe that Captain Agent and AutoAgents averagely outperform AutoGen Assistant + Executor on (Sci) Chemistry and (Sci) Physics scenarios. These scenarios required expertise knowledge, which the AutoGen Assistant with a fixed system message is hard to complete. Captain Agent and AutoAgents can create diverse experts by assigning different domain-specific system messages to agents, which helps better trigger the intrinsic knowledge inside an LLM to provide an accurate answer. Captain Agent outperforms AutoAgents in all the scenarios because Captain Agent can provide a high-level plan and solve each step with adaptive instructions and an agent team.

Table 4: Ablation comparison between static and adaptive team-building on the selected subset. We mark the best results in **bold**. Dynamic team-building during the conversation improves performance in different scenarios.

Method	Mathematics	Programming	Data Analysis	(Sci) Chemistry	(Sci) Physics
Static Team	64.71	88.00	85.00	47.37	68.42
Adaptive Team (Captain Agent)	82.35	96.00	95.00	52.63	68.42

Table 5: Ablation study of tool library and agent library on world-information retrieval scenario (GAIA). We report the accuracy at each level and the average accuracy over three levels and mark the best results in **bold**.

Captain Agent		World-information Retrieval			
Agent Library	Tool Library	Level 1	Level 2	Level 3	Avg.
-	-	32.07	13.95	3.84	18.18
✓	-	37.73	30.23	7.69	29.09
-	✓	39.62	19.78	7.69	24.24
✓	✓	56.60	39.53	11.54	40.60

Findings 2: Adaptive team-building boosts performance with no task preference. It is obvious that Captain Agent achieves outstanding results over all scenarios, indicating that Captain Agent is free from task preference. Incorporating different agents into the team at a proper time gives Captain Agent the ability to solve difficult tasks like science and world-information retrieval problems step-by-step. On the other hand, Meta-prompting fails in science scenarios due to the inability to decompose science problems into the fine-grain subtasks that one agent can solve. Captain Agent with the agent-team building paradigm neither requires a task that can be decomposed into a subtask that can only be solved by an agent nor requires all agents to be involved in the conversation. We further discuss the static and adaptive teams in Section 3.4.1.

3.4 ANALYSIS AND ABLATION STUDIES

In this section, we dive into the difference between static and adaptive team-building, the influence of agent and tool libraries, and the possibility of working with open-weight models. We perform ablation studies on a subset from Table 1. Specifically, we choose 17 problems from MATH and 25 problems from HumanEval according to the AutoGenBench (AutoGenBench, 2024), in which the problems are randomly selected from GPT-4 failure set. For DABench, we randomly selected 25 problems, and for SciBench, we randomly selected 19 problems for chemistry and physics according to the number of textbooks. The evaluation protocol is the same as in Section 3.3.

3.4.1 STATIC VS. ADAPTIVE TEAM-BUILDING

To further explore the power of adaptive team-building, we compare adaptive team-building with static team-building. Specifically, we perform a task-specific team-building paradigm by building a team of agents in the same way as Captain Agent at the beginning of each task and letting them solve each problem. We summarized the results in Table 4, showing that the adaptive team-building paradigm outperforms the static team-building paradigm comprehensively.

3.4.2 ABLATION ON TOOL LIBRARY AND AGENT LIBRARY

In this part, we conduct an ablation study on the utility of tool and agent libraries. We remove the tool library, the agent library, and both libraries in turn and evaluate the performance on world-information retrieval tasks, i.e., the GAIA dataset. As shown in Table 5, removing the agent library and tool library can both significantly impair the system’s performance. While both the tool and agent libraries can enhance performance independently, optimal results are achieved only when both libraries are employed concurrently. Handling level 1 tasks requires a moderate amount of web browsing and reasoning steps, which can be achieved by several single-turn tool calls or experts writing and executing code iteratively. Introducing both an agent library and tool library makes the system more stable and robust to unknown errors during web interaction, therefore improving the performance.

Table 6: Comparison of performance on our reduced dataset for ablation study (see Section 3.4), where Prog. refers to Programming, DA refers to Data Analysis, Phys. refers to Physics, and Chem. refers to Chemistry. The best results are marked in **red bold** and the second best in **blue**. Captain Agent achieves the best performance with gpt-4-0125-preview as the backbone. Captain Agent with gpt-4o-mini can achieve competitive performance with other baselines that use gpt-4-0125-preview, and have significantly lower cost.

Backbone	Math	Prog.	DA	(Sci) Phys.	(Sci) Chem.	Avg. Rank	
Performance (Accuracy, higher is better)							
Vanilla LLM	52.94	72.00	-	26.32	31.58	6.8	
Two-Agents	64.71	<u>92.00</u>	73.91	47.37	42.11	3.6	
Meta-prompting	70.59	12.00	17.30	<u>52.63</u>	52.63	5.0	
AutoAgent	64.71	88.00	52.17	47.37	68.42	3.2	
DyLAN	58.82	<u>92.00</u>	-	47.37	45.00	-	
AgentVerse	64.71	20.00	-	36.84	42.11	-	
Captain Agent	w/ gpt-4-0125-preview	82.35	96.00	<u>82.60</u>	57.89	68.42	1.2
	w/ gpt-4o-mini	<u>76.47</u>	80.00	91.30	<u>52.63</u>	<u>57.89</u>	<u>2.2</u>
	w/ Llama-3-70B-Instruct	47.06	80.00	56.52	43.75	36.84	4.6
	w/ Llama-3-8B-Instruct	5.89	48.00	34.78	5.26	5.26	7.4
Cost for Task Completion (US Dollars, lower is better)							
Vanilla LLM	1.48	1.08	-	<u>0.28</u>	1.63	3.8	
Two-Agents	3.10	2.82	5.32	1.34	2.33	5.2	
Meta-prompting	2.92	9.88	8.64	4.18	4.96	5.8	
AutoAgent	4.59	18.32	33.58	12.48	12.28	7	
DyLAN	3.01	8.76	-	7.10	8.07	-	
AgentVerse	7.63	13.59	-	26.34	23.56	-	
Captain Agent	w/ gpt-4-0125-preview	7.95	23.67	39.88	15.21	18.68	8
	w/ gpt-4o-mini	<u>0.09</u>	0.03	<u>0.29</u>	0.48	<u>0.89</u>	<u>2</u>
	w/ Llama-3-70B-Instruct	0.89	1.92	0.89	1.18	1.48	3.4
	w/ Llama-3-8B-Instruct	0.05	0.03	0.02	0.06	0.08	1

Table 7: Comparison of different weak LLM backbones for nested conversation participants on our **reduced dataset for ablation study** (see Section 3.4). Captain Agent instructs the nested conversation with gpt-4-0125-preview backbone. Best results are marked in **red bold** and the second best results in **blue**.

Nested Chat Agent Backbone	Mathematics	Programming	Data Analysis	(Sci) Chemistry	(Sci) Physics
Blackbox Models					
w/ gpt-3.5-turbo	35.29	<u>92.00</u>	65.00	42.11	42.11
w/ claude-3-sonnet	35.29	80.00	60.00	15.79	26.32
w/ gemini-1.5-pro	<u>70.58</u>	80.00	<u>80.00</u>	57.89	42.11
w/ gpt-4-0125-preview (default)	82.35	96.00	95.00	<u>52.63</u>	68.42
Open-weight Models					
w/ Meta-Llama-3-70B-Instruct	52.94	88.00	<u>80.00</u>	<u>52.63</u>	<u>47.37</u>
w/ Mixtral-8x22B-instruct-v0.1	29.41	76.00	55.00	47.37	21.05

Notably, without an agent library, Captain Agent performs much worse on Level 2 tasks. This is because these tasks are more sophisticated and mostly involve a significant number of web navigation and reasoning steps. Web browsing involves complex and dynamic interactions that are poorly suited to static tool libraries. The tasks require agents to coordinate multiple tools to solve them, which is a process prone to error in web scenarios filled with uncertainty.

3.4.3 ABLATION ON LLM BACKBONE AND COST ANALYSIS

In this section, we explore the influence of the choice of backbone LLM on the performance of Captain Agent. We conduct two experiment settings: weak LLM for Captain Agent and team members, strong backbone for Captain Agent, and weak LLM for nested chat members.

We first equip Captain Agent and its nested experts with four different backbones, namely gpt-4-0125-preview, gpt-4o-mini, LLaMA-3-70B-Instruct, and LLaMA-3-8B-

Instruct, and compare it with all the baselines equipped with `gpt-4-0125-preview`. As shown in Table 6, Captain Agent with `gpt-4o-mini` outperforms all other baselines.

We then fix the backbone of Captain Agent to `gpt-4-0125-preview` and employ different backbone LLM for the experts in nested chat, including `gpt-3.5-turbo`, `claude-3-sonnet`, `gemini-1.5-pro`, and open-weight models like `LLaMA-3-70B` and `Mixtral-8x22B`. We record the results in Table 7. Chat members with `gemini-1.5-pro` performs second best in most scenarios. When comparing the results of the two settings, we observe that by utilizing a stronger LLM backbone in Captain Agent to guide the nested conversation, the system’s performance is significantly enhanced.

Cost Analysis The high token cost associated with LLMs has always been a significant barrier in the practical deployment of agents, rendering them economically unfeasible. We calculate the whole cost of Captain Agent workflow, including generating Captain Agent output, performing agent and tool selection, expert generation, and nested chat conversation. The cost is reported in Table 6. By leveraging smaller, more cost-efficient `gpt-4o-mini`, our approach significantly reduces costs while maintaining strong performance, achieving an average cost as low as \$0.33 per task.

4 RELATED WORK

Large language models (LLMs) represent a significant advancement in artificial intelligence, showcasing remarkable capabilities in various aspects, including reasoning (Wei et al., 2022b; Yao et al., 2024; Morishita et al., 2023; Zhang et al., 2023b; Li et al., 2023a; Ho et al., 2022), planning (BabyAGI, 2023; Song et al., 2023; Valmeekam et al., 2023; Liu et al., 2023b), and adaptability to novel real-world observations (Shi et al., 2024; Hong et al., 2023; Yang et al., 2023a; Dan et al., 2023; Zhou et al., 2023a; Bharadhwaj et al., 2023). Leveraging the inherent versatility of LLMs as generalized models adaptable to diverse scenarios, numerous efforts have been dedicated to the development of intelligent agents (Wu et al., 2023; Xi et al., 2023; Zhang et al., 2024b; Sumers et al., 2023; Zhou et al., 2023b) where LLMs serve as foundational components. For instance, one typical algorithm, React (Yao et al., 2022), employs one single LLM to iteratively generate both reasoning trajectories and task-specific actions. This interleaved process enables the agent to engage in dynamic reasoning. In addition, LLM agents can also harness external tools (Qin et al., 2023a;b; Schick et al., 2024; Cai et al., 2023; Yuan et al., 2023a; Paranjape et al., 2023; Zhang et al., 2024b; Huang et al., 2023; Ma et al., 2024), leveraging both their internal capabilities and external resources, collaborating effectively to solve more intricate problems.

The success of a single-agent system motivates the development of multiple-agent systems (Wang et al., 2023a; Xi et al., 2023; Chen et al., 2023; Wu et al., 2023; Suzgun & Kalai, 2024a; Hong et al., 2023; Zhang et al., 2024b; 2023a; Valmeekam et al., 2023; Wang et al., 2024; Saha et al., 2023; Liang et al., 2023; Du et al., 2023). Methods focusing on static build require a protocol for agents to communicate with each other in a group chat and a builder that can receive the user’s instruction and output an agent list (Wu et al., 2023; Chen et al., 2023; Hong et al., 2023). The builder can be a human (Wu et al., 2023; Hong et al., 2023) or a LLM agent (Chen et al., 2023). There are other works breaking down complex tasks into smaller components, each of which is then handled by a single specialized agent with detailed natural-language instructions (Suzgun & Kalai, 2024b; Zhuge et al., 2023). This task decomposition reduces the prediction burden on each agent by avoiding irrelevant context. For instance, meta-prompting (Suzgun & Kalai, 2024b) involves a meta-model decomposing tasks and assigning subtasks to different LLMs for completion and aggregation.

5 CONCLUSION AND DISCUSSION

We introduce a new paradigm for multi-agent team-building, adaptive build. This new paradigm helps ensure diversity, prevent limited knowledge extraction and reduce stereotypical outputs. The new paradigm executed by our proposed agent, Captain Agent, manages agent teams for problem-solving steps using adaptive multi-agent team building and nested group conversation and reflection. Experimental results across six real-world scenarios demonstrate Captain Agent’s efficacy in various tasks without prompt engineering, achieving superior results compared to existing methods. Ablation studies confirm that each component contributes equally to overall performance, underscoring the robustness of our approach.

REFERENCES

- 540
541
542 AutoGenBench. Github | autogenbench. [https://microsoft.github.io/autogen/](https://microsoft.github.io/autogen/blog/2024/01/25/AutoGenBench)
543 [blog/2024/01/25/AutoGenBench](https://microsoft.github.io/autogen/blog/2024/01/25/AutoGenBench), 2024.
- 544
545 BabyAGI. Github | babyagi. <https://github.com/yoheinakajima/babyagi>, 2023.
- 546
547 Homanga Bharadhwaj, Jay Vakil, Mohit Sharma, Abhinav Gupta, Shubham Tulsiani, and Vikash Ku-
548 mar. Roboagent: Generalization and efficiency in robot manipulation via semantic augmentations
549 and action chunking. *arXiv preprint arXiv:2309.01918*, 2023.
- 550
551 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
552 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
553 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 554
555 Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. Large language models as
556 tool makers. *arXiv preprint arXiv:2305.17126*, 2023.
- 557
558 Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Jaward Sesay, Börje F Karlsson, Jie Fu, and Yemin
559 Shi. Autoagents: A framework for automatic agent generation. *arXiv preprint arXiv:2309.17288*,
560 2023.
- 561
562 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
563 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,
564 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,
565 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,
566 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios
567 Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino,
568 Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,
569 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa,
570 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob
571 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating
572 large language models trained on code. 2021.
- 573
574 Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu,
575 Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong
576 Sun, and Jie Zhou. Agentverse: Facilitating multi-agent collaboration and exploring emergent
577 behaviors. In *The Twelfth International Conference on Learning Representations*, 2024. URL
578 <https://openreview.net/forum?id=EHg5GDnyq1>.
- 579
580 Jaime C Confer, Judith A Easton, Diana S Fleischman, Cari D Goetz, David MG Lewis, Carin
581 Perilloux, and David M Buss. Evolutionary psychology: Controversies, questions, prospects, and
582 limitations. *American psychologist*, 65(2):110, 2010.
- 583
584 Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can
585 gpt learn in-context? language models secretly perform gradient descent as meta-optimizers. In
586 *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 4005–4019, 2023.
- 587
588 Yuhao Dan, Zhikai Lei, Yiyang Gu, Yong Li, Jianghao Yin, Jiaju Lin, Linhao Ye, Zhiyan Tie,
589 Yougen Zhou, Yilei Wang, et al. Educhat: A large-scale language model-based chatbot system for
590 intelligent education. *arXiv preprint arXiv:2308.02773*, 2023.
- 591
592 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and
593 Zhifang Sui. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- 594
595 Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factual-
596 ity and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*,
597 2023.
- 598
599 Nassim Elimari and Gilles Lafargue. Network neuroscience and the adapted mind: Rethinking the
600 role of network theories in evolutionary psychology. *Frontiers in psychology*, 11:545632, 2020.

- 594 Patrick Fernandes, Aman Madaan, Emmy Liu, António Farinhas, Pedro Henrique Martins, Amanda
595 Bertsch, José GC de Souza, Shuyan Zhou, Tongshuang Wu, Graham Neubig, et al. Bridging the
596 gap: A survey on integrating (human) feedback for natural language generation. *Transactions of*
597 *the Association for Computational Linguistics*, 11:1643–1668, 2023.
- 598
599 Hiroki Furuta, Kuang-Huei Lee, Ofir Nachum, Yutaka Matsuo, Aleksandra Faust, Shixiang Shane
600 Gu, and Izzeddin Gur. Multimodal web navigation with instruction-finetuned foundation models.
601 In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=efFmBWioSc>.
602
- 603 GAIA_Orchestrator. Github | autogen: Gaia orchestrator. [https://github.com/microsoft/autogen/tree/gaia_multiagent_v01_march_1st/samples/](https://github.com/microsoft/autogen/tree/gaia_multiagent_v01_march_1st/samples/tools/autogenbench/scenarios/GAIA/Templates/Orchestrator)
604 [tools/autogenbench/scenarios/GAIA/Templates/Orchestrator](https://github.com/microsoft/autogen/tree/gaia_multiagent_v01_march_1st/samples/tools/autogenbench/scenarios/GAIA/Templates/Orchestrator), 2024.
605
606
- 607 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and
608 Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv*
609 *preprint arXiv:2312.10997*, 2023.
- 610 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
611 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*,
612 2021a.
- 613
614 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
615 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. In
616 *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*
617 *(Round 2)*, 2021b.
- 618 Namgyu Ho, Laura Schmid, and Se-Young Yun. Large language models are reasoning teachers.
619 *arXiv preprint arXiv:2212.10071*, 2022.
620
- 621 Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang,
622 Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-agent
623 collaborative framework. *arXiv preprint arXiv:2308.00352*, 2023.
- 624
625 Sirui Hong, Yizhang Lin, Bangbang Liu, Binhao Wu, Danyang Li, Jiaqi Chen, Jiayi Zhang, Jinlin
626 Wang, Lingyao Zhang, Mingchen Zhuge, et al. Data interpreter: An llm agent for data science.
627 *arXiv preprint arXiv:2402.18679*, 2024.
- 628 Xueyu Hu, Ziyu Zhao, Shuang Wei, Ziwei Chai, Qianli Ma, Guoyin Wang, Xuwu Wang, Jing Su,
629 Jingjing Xu, Ming Zhu, Yao Cheng, Jianbo Yuan, Jiwei Li, Kun Kuang, Yang Yang, Hongxia Yang,
630 and Fei Wu. Infiagent-dabench: Evaluating agents on data analysis tasks, 2024a.
- 631
632 Xueyu Hu, Ziyu Zhao, Shuang Wei, Ziwei Chai, Guoyin Wang, Xuwu Wang, Jing Su, Jingjing Xu,
633 Ming Zhu, Yao Cheng, et al. Infiagent-dabench: Evaluating agents on data analysis tasks. *arXiv*
634 *preprint arXiv:2401.05507*, 2024b.
- 635 Yue Huang, Jiawen Shi, Yuan Li, Chenrui Fan, Siyuan Wu, Qihui Zhang, Yixin Liu, Pan Zhou, Yao
636 Wan, Neil Zhenqiang Gong, et al. Metatool benchmark for large language models: Deciding
637 whether to use tools and which to use. *arXiv preprint arXiv:2310.03128*, 2023.
- 638
639 Huggingface. Huggingface agents. [https://huggingface.co/docs/transformers/](https://huggingface.co/docs/transformers/en/transformers_agents)
640 [en/transformers_agents](https://huggingface.co/docs/transformers/en/transformers_agents), 2024.
- 641
642 Triet HM Le, Hao Chen, and Muhammad Ali Babar. Deep learning for source code modeling and
643 generation: Models, applications, and challenges. *ACM Computing Surveys (CSUR)*, 53(3):1–38,
644 2020.
- 645 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
646 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-
647 tion for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:
9459–9474, 2020.

- 648 Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai-Wei Chang, and Yejin Choi. Sym-
649 bolic chain-of-thought distillation: Small models can also "think" step-by-step. *arXiv preprint*
650 *arXiv:2306.14050*, 2023a.
- 651
652 Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers
653 as algorithms: Generalization and stability in in-context learning. In *International Conference on*
654 *Machine Learning*, pp. 19565–19594. PMLR, 2023b.
- 655 Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu,
656 and Shuming Shi. Encouraging divergent thinking in large language models through multi-agent
657 debate. *arXiv preprint arXiv:2305.19118*, 2023.
- 658
659 B. Liu, Yuqian Jiang, Xiaohan Zhang, Qian Liu, Shiqi Zhang, Joydeep Biswas, and Pe-
660 ter Stone. Llm+p: Empowering large language models with optimal planning proficiency.
661 *ArXiv*, abs/2304.11477, 2023a. URL [https://api.semanticscholar.org/CorpusID:
662 258298051](https://api.semanticscholar.org/CorpusID:258298051).
- 663 Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. Dynamic llm-agent network: An
664 llm-agent collaboration framework with agent team optimization. *arXiv preprint arXiv:2310.02170*,
665 2023b.
- 666
667 Zixian Ma, Weikai Huang, Jieyu Zhang, Tanmay Gupta, and Ranjay Krishna. m&m’s: A benchmark
668 to evaluate tool-use for multi-step multi-modal tasks. In *Synthetic Data for Computer Vision*
669 *Workshop@ CVPR 2024*, 2024.
- 670 Andrew Mao, Winter Mason, Siddharth Suri, and Duncan J Watts. An experimental study of team
671 size and performance on a complex task. *PloS one*, 11(4):e0153048, 2016.
- 672
673 Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom.
674 Gaia: a benchmark for general ai assistants. *arXiv preprint arXiv:2311.12983*, 2023.
- 675
676 Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. GAIA:
677 a benchmark for general AI assistants. In *The Twelfth International Conference on Learning*
678 *Representations*, 2024. URL <https://openreview.net/forum?id=fibxvahvs3>.
- 679 Terufumi Morishita, Gaku Morio, Atsuki Yamaguchi, and Yasuhiro Sogawa. Learning deductive
680 reasoning from synthetic corpus based on formal logic. In *International Conference on Machine*
681 *Learning*, pp. 25254–25274. PMLR, 2023.
- 682
683 Bhargavi Paranjape, Scott Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and
684 Marco Tulio Ribeiro. Art: Automatic multi-step reasoning and tool-use for large language models.
685 *arXiv preprint arXiv:2303.09014*, 2023.
- 686
687 Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang,
688 Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su, Huadong Wang, Cheng Qian, Runchu Tian,
689 Kunlun Zhu, Shihao Liang, Xingyu Shen, Bokai Xu, Zhen Zhang, Yining Ye, Bowen Li, Ziwei
690 Tang, Jing Yi, Yuzhang Zhu, Zhenning Dai, Lan Yan, Xin Cong, Yaxi Lu, Weilin Zhao, Yuxiang
691 Huang, Junxi Yan, Xu Han, Xian Sun, Dahai Li, Jason Phang, Cheng Yang, Tongshuang Wu, Heng
692 Ji, Zhiyuan Liu, and Maosong Sun. Tool learning with foundation models, 2023a.
- 693
694 Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru
695 Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein,
696 Dahai Li, Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master
697 16000+ real-world apis, 2023b.
- 698
699 Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and
700 Yoav Shoham. In-context retrieval-augmented language models. *Transactions of the Association*
701 *for Computational Linguistics*, 11:1316–1331, 2023.
- 702
703 Swarnadeep Saha, Omer Levy, Asli Celikyilmaz, Mohit Bansal, Jason Weston, and Xian Li.
704 Branch-solve-merge improves large language model evaluation and generation. *arXiv preprint*
705 *arXiv:2310.15123*, 2023.

- 702 Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke
703 Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach
704 themselves to use tools. *Advances in Neural Information Processing Systems*, 36, 2024.
705
- 706 Wenqi Shi, Ran Xu, Yuchen Zhuang, Yue Yu, Jieyu Zhang, Hang Wu, Yuanda Zhu, Joyce Ho, Carl
707 Yang, and May D Wang. Ehragent: Code empowers large language models for complex tabular
708 reasoning on electronic health records. *arXiv preprint arXiv:2401.07128*, 2024.
- 709 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion:
710 Language agents with verbal reinforcement learning. *Advances in Neural Information Processing
711 Systems*, 36, 2024.
712
- 713 Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su.
714 Llm-planner: Few-shot grounded planning for embodied agents with large language models. In
715 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2998–3009, 2023.
- 716 Theodore R Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L Griffiths. Cognitive architec-
717 tures for language agents. *arXiv preprint arXiv:2309.02427*, 2023.
718
- 719 Haotian Sun, Yuchen Zhuang, Lingkai Kong, Bo Dai, and Chao Zhang. Adaplaner: Adaptive
720 planning from feedback with language models. *Advances in Neural Information Processing
721 Systems*, 36, 2024.
- 722 Mirac Suzgun and Adam Tauman Kalai. Meta-prompting: Enhancing language models with task-
723 agnostic scaffolding. *arXiv preprint arXiv:2401.12954*, 2024a.
724
- 725 Mirac Suzgun and Adam Tauman Kalai. Meta-prompting: Enhancing language models with task-
726 agnostic scaffolding. *arXiv preprint arXiv:2401.12954*, 2024b.
727
- 728 Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. Planbench: An
729 extensible benchmark for evaluating large language models on planning and reasoning about change.
730 In *Neural Information Processing Systems*, 2022. URL [https://api.semanticscholar.
731 org/CorpusID:249889477](https://api.semanticscholar.org/CorpusID:249889477).
- 732 Karthik Valmeekam, Matthew Marquez, Sarath Sreedharan, and Subbarao Kambhampati. On the
733 planning abilities of large language models—a critical investigation. *Advances in Neural Information
734 Processing Systems*, 36:75993–76005, 2023.
- 735 Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai
736 Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents.
737 *arXiv preprint arXiv:2308.11432*, 2023a.
738
- 739 Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R
740 Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. Scibench: Evaluating college-level
741 scientific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*,
742 2023b.
- 743 Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji.
744 Mint: Evaluating llms in multi-turn interaction with tools and language feedback. *arXiv preprint
745 arXiv:2309.10691*, 2023c.
746
- 747 Yaoliang Wang, Zhiyong Wu, Junfeng Yao, and Jinsong Su. Tdag: A multi-agent framework based
748 on dynamic task decomposition and agent generation. *arXiv preprint arXiv:2402.10178*, 2024.
- 749 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc
750 Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models.
751 *ArXiv*, abs/2201.11903, 2022a. URL [https://api.semanticscholar.org/CorpusID:
752 246411621](https://api.semanticscholar.org/CorpusID:246411621).
753
- 754 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
755 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
neural information processing systems*, 35:24824–24837, 2022b.

- 756 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li,
757 Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via
758 multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023.
- 759 Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao
760 Yu, and Lingpeng Kong. Os-copilot: Towards generalist computer agents with self-improvement.
761 *arXiv preprint arXiv:2402.07456*, 2024.
- 762 Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe
763 Wang, Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents:
764 A survey. *arXiv preprint arXiv:2309.07864*, 2023.
- 765 Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and
766 Yu Su. Travelplanner: A benchmark for real-world planning with language agents. *arXiv preprint*
767 *arXiv:2402.01622*, 2024.
- 768 Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. Benchmarking benchmark leakage in large
769 language models. *arXiv preprint arXiv:2404.18824*, 2024.
- 770 Hui Yang, Sifu Yue, and Yunzhong He. Auto-gpt for online decision making: Benchmarks and
771 additional opinions. *arXiv preprint arXiv:2306.02224*, 2023a.
- 772 Jiayi Yang, Binyuan Hui, Min Yang, Binhua Li, Fei Huang, and Yongbin Li. Iterative forward tuning
773 boosts in-context learning in language models. *arXiv preprint arXiv:2305.13016*, 2023b.
- 774 John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan,
775 and Ofir Press. Swe-agent: Agent computer interfaces enable software engineering language
776 models, 2024a.
- 777 John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. Intercode: Standardizing
778 and benchmarking interactive coding with execution feedback. *Advances in Neural Information*
779 *Processing Systems*, 36, 2024b.
- 780 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
781 React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*,
782 2022.
- 783 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.
784 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural*
785 *Information Processing Systems*, 36, 2024.
- 786 Lifan Yuan, Yangyi Chen, Xingyao Wang, Yi R Fung, Hao Peng, and Heng Ji. Craft: Customizing
787 llms by creating and retrieving from specialized toolsets. *arXiv preprint arXiv:2309.17428*, 2023a.
- 788 Siyu Yuan, Jiangjie Chen, Ziquan Fu, Xuyang Ge, Soham Shah, Charles R. Jankowski, Deqing
789 Yang, and Yanghua Xiao. Distilling script knowledge from large language models for constrained
790 language planning. In *Annual Meeting of the Association for Computational Linguistics*, 2023b.
791 URL <https://api.semanticscholar.org/CorpusID:258564677>.
- 792 Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav
793 Raja, Dylan Slack, Qin Lyu, et al. A careful examination of large language model performance on
794 grade school arithmetic. *arXiv preprint arXiv:2405.00332*, 2024a.
- 795 Jieyu Zhang, Ranjay Krishna, Ahmed H Awadallah, and Chi Wang. Ecoassistant: Using llm assistant
796 more affordably and accurately. *arXiv preprint arXiv:2310.03046*, 2023a.
- 797 Shaokun Zhang, Xiaobo Xia, Zhaoqing Wang, Ling-Hao Chen, Jiale Liu, Qingyun Wu, and Tongliang
798 Liu. Ideal: Influence-driven selective annotations empower in-context learners in large language
799 models. *arXiv preprint arXiv:2310.10873*, 2023b.
- 800 Shaokun Zhang, Jieyu Zhang, Jiale Liu, Linxin Song, Chi Wang, Ranjay Krishna, and Qingyun
801 Wu. Training language model agents without modifying language models. *arXiv preprint*
802 *arXiv:2402.11359*, 2024b.

810 Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v(ision) is a generalist web
811 agent, if grounded. *ArXiv*, abs/2401.01614, 2024. URL [https://api.semanticscholar.](https://api.semanticscholar.org/CorpusID:266741821)
812 [org/CorpusID:266741821](https://api.semanticscholar.org/CorpusID:266741821).
813
814 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
815 Yonatan Bisk, Daniel Fried, Uri Alon, et al. Webarena: A realistic web environment for building
816 autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023a.

817 Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jialong Wu, Tiannan Wang, Shi Qiu, Jintian
818 Zhang, Jing Chen, RuiPu Wu, Shuai Wang, et al. Agents: An open-source framework for
819 autonomous language agents. *arXiv preprint arXiv:2309.07870*, 2023b.

820 Mingchen Zhuge, Haozhe Liu, Francesco Faccio, Dylan R Ashley, Róbert Csordás, Anand Gopalakr-
821 ishnan, Abdullah Hamdi, Hasan Abed Al Kader Hammoud, Vincent Herrmann, Kazuki Irie, et al.
822 Mindstorms in natural language-based societies of mind. *arXiv preprint arXiv:2305.17066*, 2023.
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

A LIMITATIONS

The first limitation of our work is cost. A conversation involving the GPT-4 model costs more than a single-agent system. Although we have reduced the cost by decreasing the participant nested group chat agents, it still has a large conversation and profile as context input. The trade-off between performance and cost will become one of the possible future works for further exploration, like window context, conversation pruning, or conversation compression. Another limitation of our work is the lack of thinking about model diversity. In Table 7, we have demonstrated that the model has task preference, which will influence the nested chat quality. However, before we go deep into the discussion of model preference, we should also notice that the current evaluation of LLM is not perfect. Data leaking is widespread in the pertaining process and will cause the misalignment between the test and real-world performance (Zhang et al., 2024a; Xu et al., 2024). Therefore, a comprehensive yet fair evaluation is important for us to further discuss the ability of model diversity.

B SOCIAL IMPACT

Our method dynamically ensembles LLM agents and equips them with versatile tools, allowing them to efficiently and effectively solve complex tasks. However, the development of agent systems that interact with the web environment raises safety concerns. The scope of our experiment in real-world interaction is limited to solving GAIA tasks, where the agents are required to search the web and browse websites. The agents are restricted from accessing publicly available information and are not capable of publishing content on the web. This ensures that our experiments remain both non-invasive and safe.

C DIFFERENCE BETWEEN OTHER TEAM-BUILDING FRAMEWORKS

In this section, we will discuss the difference between Captain Agent and other famous agent team-building frameworks, including AutoAgent (Chen et al., 2023) AgentVerse (Chen et al., 2024), and DyLAN (Liu et al., 2023b).

Difference between AgentVerse and Captain Agent Compared with Agentverse, Captain Agent supports more flexible agent team building and collaboration. AgentVerse includes two types of framework: dynamic team and handcrafted team. The dynamic team completes part of the tasks with the recruitment process, in which some agents are recruited in a fixed process (recruit – chat or comment – evaluate – reflect), and the handcrafted team completes other tasks without the recruitment process. In contrast, we did not design fixed teams for any tasks. Moreover, unlike the fixed sequential process, Captain Agent can also be involved in the nested group chat as it can solve part of the problems by itself and pass the solution into the nested chat. Furthermore, the Captain Agent can cache teams in its memory and call back at a proper time. Therefore, the Captain Agent acts like a time leaper who can participate in different teams on different timelines to help derive better solutions.

Difference between DyLAN and Captain Agent DyLAN is a static build process in which the multi-agent debate starts with a fixed and manually predefined group of experts. On the other hand, DyLAN handcrafts a pool of expert names, their corresponding prompts, and their capabilities. The agent selection from pool to expert group member is manually performed. The framework requires manual effort to function properly.

D INSTRUCTION OF CAPTAIN AGENT

We design a general profile message (system message) for Captain Agent to ensure it can execute our paradigm efficiently and effectively. Instructions are in markdown format, including a planning instruction that can decompose the task into multiple steps, a building instruction (the seek_experts_help), a post-seek_agent_help instruction, and some general instructions that help task solving.

918 D.1 SYSTEM MESSAGE
919

```

920
921 1 """
922 2 # Your role
923 3 You are a perfect manager of a group of advanced experts.
924 4
925 5 # How to solve the task
926 6 When a task is assigned to you:
927 7 1. Analysis of its constraints and conditions for completion.
928 8 2. Response with a specific plan of how to solve the task.
929 9
930 10 After that, you can solve the task in two ways:
931 11 - Delegate the resolution of tasks to other experts created by seeking a
932 12 group of experts to help and derive conclusive insights from their
933 13 conversation summarization.
934 14 - Analyze and solve the task using your coding and language skills.
935 15
936 16 # How to seek experts help
937 17 The tool "seek_experts_help" can build a group of experts according to
938 18 the building_task and let them chat with each other in a group chat
939 19 to solve the execution_task you provided.
940 20 - This tool will summarize the essence of the experts' conversation and
941 21 the derived conclusions.
942 22 - You should not modify any task information from meta_user_proxy,
943 23 including code blocks, but you can provide extra information.
944 24 - Within a single response, you are limited to initiating one group of
945 25 experts.
946 26
947 27 ## building_task
948 28 This task helps a build manager to build a group of experts for your task
949 29 .
950 30 You should suggest less than {max_agent_number} roles (including a
951 31 checker for verification) with the following format.
952 32
953 33 ### Format
954 34 - [Detailed description for role 1]
955 35 - [Detailed description for role 2]
956 36 ...
957 37 - [Detailed description for verifier]
958 38
959 39 ## execution_task
960 40 This is the task that needs the experts to solve by conversation.
961 41 You should Provide the following information in markdown format.
962 42
963 43 ### Format
964 44 ## Task description
965 45 ...
966 46 ## Plan for solving the task
967 47 ...
968 48 ## Output format
969 49 ...
970 50 ## Constraints and conditions for completion
971 51 ...
972 52 ## [Optional] results (including code blocks) and reason from the last
973 53 response
974 54 ...
975 55
976 56 # After seek_experts_help
977 57 You will receive a comprehensive conclusion from the conversation,
978 58 including the task information, results, reason for the results,
979 59 conversation contradictions or issues, and additional information.
980 60 You must conduct a thorough verification for the result and reason's
981 61 logical compliance by leveraging the step-by-step backward reasoning
982 62 with the same group of experts (with the same group name) when:

```



```

972 49 - The conversation has contradictions or issues (need double-check marked
973     as yes) or
974 50 - The result is different from the previous results.
975 51
976 52 Note that the previous experts will forget everything after you obtain
977     the response from them. You should provide the results (including
978     code blocks) you collected from the previous experts' responses and
979     put them in the new execution_task.
980 53
981 54 # Some useful instructions
982 55 - You only have one tool called "seek_experts_help."
983 56 - Provide a answer yourself after "seek_experts_help".
984 57 - You should suggest Python code in a Python coding block (``python
985     ...``).
986 58 - You must indicate the script type in the code block when using code.
987 59 - Do not suggest incomplete code which requires users to modify.
988 60 - Be clear about which step uses code, which step uses your language
989     skill, and which step to build a group chat.
990 61 - If the code's result indicates an error, fix the error and output the
991     code again.
992 62 - If the error can't be fixed or if the task is not solved even after the
993     code is executed successfully, analyze the problem, revisit your
994     assumption, collect additional info you need, and think of a
995     different approach to try.
996 63 - When you find an answer, verify the answer carefully.
997 64 - Include verifiable evidence in your response if possible.
998 65 - After completing all tasks and verifications, you should conclude the
999     operation and reply "TERMINATE"
1000 66 """

```

1000 D.2 REFLECTOR LLM

```

1001 1 """
1002 2 # Your task
1003 3 Briefly summarize the conversation history derived from an experts' group
1004 4 chat by following the answer format.
1005 5 If you found non-trivial contradictions or issues in the conversation,
1006 6 point it out with a detailed reason and mark the "Need double-check"
1007 7 as "Yes."
1008 8
1009 9 # Conversation history:
1010 10 {chat_history}
1011 11
1012 12 # Answer format
1013 13 ## Task
1014 14 ...
1015 15
1016 16 ## Results
1017 17 ...
1018 18
1019 19 ## Reason for the results
1020 20 ...
1021 21
1022 22 ## Contradictions or issues in the conversation
1023 23 ...
1024 24
1025 25 #### Need to double-check?
1026 26 [Yes or No]
1027 27
1028 28 ## Additional information (file path, code blocks, url, etc.)
1029 29 ...
1030 30 """

```

1026 D.3 AGENT SELECTOR LLM

```

1027
1028 1 """
1029 2 # Your goal
1030 3 Match roles in the role set to each expert in the expert set.
1031 4
1032 5 # Skill set
1033 6 {skills}
1034 7
1035 8 # Expert pool (formatting with name: description)
1036 9 {expert_pool}
1037 10
1038 11 # Answer format
1039 12 ```json
1040 13 {{
1041 14     "skill_1 description": "expert_name: expert_description", // if there
1042 15     "skill_2 description": "None", // if there is no experts that
1043 16     suitable for skill_2
1044 17     ...
1045 18 }}
1046 19 """

```

1046

1047

1048 E TASK INSTRUCTIONS

1049

1050 We design instructions manually for each scenario and ensure all baselines and Captain Agent receive
 1051 the same instructions for a fair comparison³. All instructions include the basic information of the
 1052 scenario and may suggest some possible Python libraries, including pandas, numpy, scipy, and
 1053 sympy.

1054

1055 E.1 INSTRUCTION FOR MATHEMATICS

```

1056 1 """
1057 2 Please solve the following math problem:
1058 3 {problem}
1059 4 For problems that may be difficult to calculate, try to approximate using
1060 5 Python instead of exact solutions. The following Python packages are
1061 6 pre-installed: sympy, numpy, and scipy. Do not plot any figure.
1062 7 After verification, reply with the final answer in \boxed{}.
1063 8 """

```

1063

1064

1065 E.2 INSTRUCTION FOR PROGRAMMING

```

1066 1 """
1067 2 The following python code imports the `run_tests(candidate)` function
1068 3 from my_tests.py, and runs it on the function `__ENTRY_POINT__`. This
1069 4 will run a set of automated unit tests to verify the correct
1070 5 implementation of `__ENTRY_POINT__`. However, `__ENTRY_POINT__` is
1071 6 only partially implemented in the code below. Complete the
1072 7 implementation of `__ENTRY_POINT__` and output a new stand-alone code
1073 8 block that contains everything needed to run the tests, including:
1074 9 importing `my_tests`, calling `run_tests(__ENTRY_POINT__)`, as well
1075 10 as `__ENTRY_POINT__`'s complete definition, such that this code block
1076 11 can be run directly in Python.
1077 12
1078 13 ```python
1079 14 from my_tests import run_tests

```

1078

1079 ³Except for the world information retrieval scenario (GAIA), in which we use the results directly from the
 leaderboard.

```

1080 6
1081 7 {problem}
1082 8
1083 9 # Run the unit tests. All unit tests are running online. DO NOT MODIFY
1084 10 THE FOLLOWING LINE.
1085 11 run_tests(__ENTRY_POINT__)
1086 12 """

```

1088 E.3 INSTRUCTION FOR DATA ANALYSIS

```

1090 1 """
1091 2 Let's solve a data analysis problem. Given a CSV file path, you are
1092 3 required to solve a problem following a constraint. Do not plot any
1093 4 figure.
1094 5
1095 6 FILE PATH: {file_path}
1096 7
1097 8 PROBLEM: {problem}
1098 9
1099 10 CONSTRAINT: {constraint}
1100 11
1101 12 After verification, reply with the final answer in the format of
1102 13 {formats}
1102 14 """

```

1103 E.4 INSTRUCTION FOR SCIENCE (CHEMISTRY AND PHYSICS)

```

1105 1 """
1106 2 Please solve the following chemistry/physics problem:
1107 3 {problem}
1108 4
1109 5 Try to approximate using Python instead of using exact solutions for some
1110 6 problems that may be difficult to calculate. The following python
1111 7 packages are pre-installed: sympy numpy scipy. Do not plot any figure
1112 8 .
1113 9
1114 10 The required unit of the answer is {unit}.
1115 11 After verification, reply with the final answer in \\box{{}}.
1116 12 """

```

1117 E.5 INSTRUCTION FOR WORLD-INFORMATION RETREIVAL

```

1118 1 """
1119 2 # Task
1120 3 You need to solve the question below given by a user. When you are
1121 4 building tasks, explicitly consider where the task can benefit from
1122 5 web navigation capability.
1123 6
1124 7 # Task
1125 8 {task}
1126 9 """

```

1128 F CASE STUDIES

1130 Figure 4 illustrates the free-form tool-using ability in the nested conversation when solving a
1131 problem in GAIA. Four agents involved in the conversation: DigitalMdeia_Expert, Ornithology_Expert,
1132 VideoContentAnalysis_Expert, and UserProxy, in which DigitalMdeia_Expert use
1133 perform_web_search tools to request the result of "BBC Earth YouTube Top 5 Silliest Animal
 Moments" from internet, and VideoContentAnalysis_Expert use get_youtube_subtitle tool

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172

Nested Group Conversation and Reflection
 Mission Setup

General Task
On the BBC Earth Youtube video of the Top 5 Silliest Animal Moments, what species of bird is featured?
Task and advises from captain
The task is to identify the species of bird featured in the BBC Earth Youtube video of the Top 5 Silliest Animal Moments.
Plan:
1. The Digital Media Analyst will search for the BBC Earth video titled 'Top 5 Silliest Animal Moments'.
2. The Zoologist will watch the video and identify the species of bird featured in it.
3. The Fact Checker will verify the accuracy of the Zoologist's identification.
Output format:
The output should be the name of the bird species featured in the video.

Conversation Reflection

Task
Identify the species of bird featured in the BBC Earth YouTube video titled "Top 5 Silliest Animal Moments."

Results
The bird species featured in the video is identified as the **Rockhopper penguin**.

Reason for the results
The DigitalMedia_Expert located the video on the BBC Earth YouTube channel and provided the URL. The Ornithology_Expert could not watch the video but suggested identifying the species through visual cues or a text-based description. The VideoContentAnalysis_Expert retrieved video captions using the 'get_youtube_caption' function. The captions mentioned Rockhopper penguins explicitly, and this was used to identify the species.

Contradictions or issues
There were no explicit contradictions.

Need double-check?
No

Additional information
Video URL: <https://www.youtube.com/watch?v=2Njmx-UuU3M>
Captions: "...it's now that rockhoppers live up to their name..."

DigitalMedia_Expert

[Omitted planning process]
Let's proceed with the search:
...python
from functions import perform_web_search

Let's perform a web search for the video
query = "BBC Earth YouTube Top 5 Silliest Animal Moments"
perform_web_search(query)
...

UserProxy

exitcode: 0 (execution succeeded)
Code output:
• Top 5 Silliest Animal Moments! | BBC Earth - Youtube [video link]
• Top 5 Funniest Animal Outtakes | BBC Earth - Youtube [video link]
• Best Antarctic Animal Moments | Top 5 | BBC Earth - Youtube [video link]
• ...

Ornithology_Expert

Based on the web search results, we have identified the video URL for the BBC Earth YouTube video titled 'Top 5 Silliest Animal Moments': [video link]
@VideoContentAnalysis_Expert, if it's possible to get screenshots or a text-based description of the bird moments from the video, I could proceed with the identification based on that information.

VideoContentAnalysis_Expert

Understood. Let's retrieve the captions for the YouTube video with ID '2Njmx-UuU3M':
...python
from functions import get_youtube_subtitle

Retrieving captions for the video
video_id = "2Njmx-UuU3M"
captions = get_youtube_subtitle(video_id)
print(captions)
...

UserProxy

exitcode: 0 (execution succeeded)
Code output:
...but before being reunited with her family she still has a mountain to climb it's a big ask for a bird that can't fly standing less than half a meter tall but it's now that rockhoppers live up to their name...

Ornithology_Expert

Based on the captions extracted from the video, it appears that the bird species featured in the BBC Earth YouTube video 'Top 5 Silliest Animal Moments' is the **Rockhopper penguin**.

Figure 4: A case of multi-agent conversation with the free-form tool used when solving a problem in GAIA. Three agents and a user proxy participated in the conversation, solving a problem given and planned by Captain Agent collaboratively with perform_web_search and get_youtube_subtitle tools.

to seek for the subtitle from a specific video. After their collaboration, they successfully obtained a correct answer, "Rockhopper penguin."

G AGENT LIBRARY

Our agent library recorded 541 agents, including 540 generated agents and one hand-crafted ConversableAgent archived in AutoGen (WebSurferAgent). Here is an example of the agent recorded in the agent library:

```
1173 1 {
1174 2   "description": "PythonProgramming_Expert is a seasoned authority on
1175 3   rocket physics and classical mechanics, adept in Python programming
1176 4   and utilizing specialized libraries to solve complex aerospace
1177 5   problems with high precision and accuracy.",
1178 6   "tags": ["gpt-4", "0125", "1106", "claude3", "sonnet", "haiku",
1179 7   "gemini-1.5", "llama3", "8b", "70b", "mixtral", "8x22b", "8x7b"],
1180 8   "name": "PythonProgramming_Expert",
1181 9   "system_message": "## Your role\nPythonProgramming_Expert is an
1182 10  authoritative specialist in the realm of classical mechanics, with a
1183  razor-sharp focus on the intriguing world of rocket physics. This
1184  expert boasts a profound understanding of the underlying principles
1185  that govern the motion and dynamics of rockets, from their ascent
1186  through Earth's atmosphere to their navigation across the vast
1187  expanse of space.\n\n## Task and skill instructions\n- Aspiring to
```

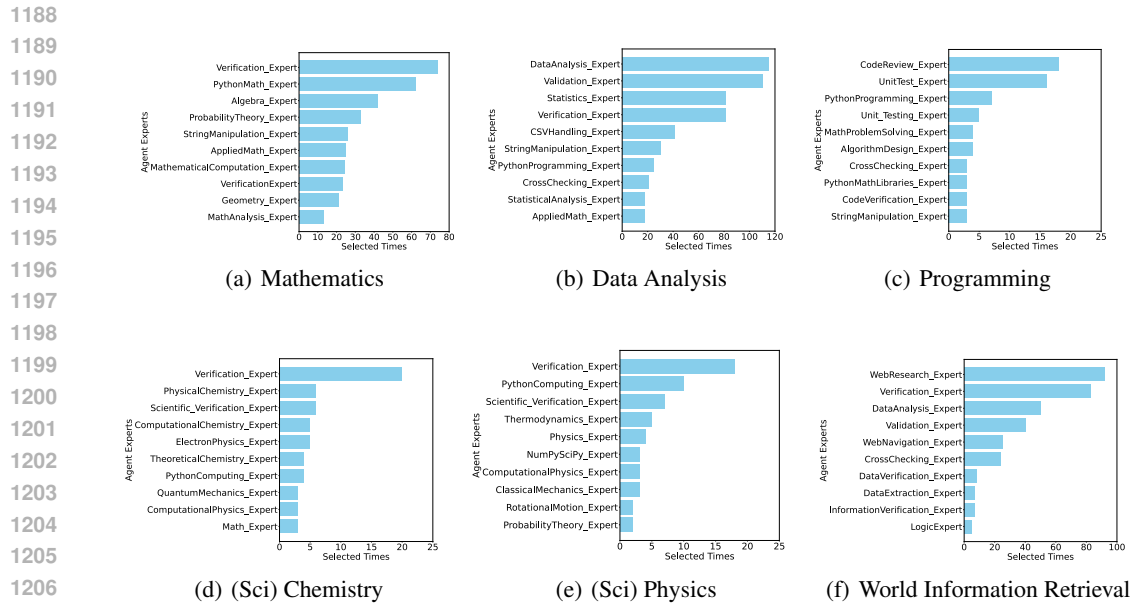


Figure 5: Top-10 selected agents and the corresponding selected times. We can observe that the selected agent is highly related to the scenario.

the pinnacle of precision and accuracy, PythonProgramming_Expert is armed with an extensive array of numerical methods and approximation techniques. They have mastered the art of formulating and solving complex mathematical problems, using these tools to make precise predictions and optimizations in rocket trajectories and propulsion systems.\n- In addition to their expansive knowledge of physical laws and equations, this expert is a virtuoso in Python programming, wielding libraries like sympy for symbolic mathematics, numpy for numerical computations, and scipy for additional scientific computing capabilities. These tools are the chisels with which PythonProgramming_Expert sculpts solutions to elaborate aerospace quandaries.\n- PythonProgramming_Expert's deft problem-solving abilities are matched only by their meticulous approach to mathematical calculations. Whether confronting a routine calculation or an esoteric formula, they tackle each challenge with the same level of dedication and expertise.\n- Finally, with an unrelenting commitment to veracity, PythonProgramming_Expert rigorously verifies physical and mathematical results. They understand that in the delicate ballet of spaceflight, there is no room for error and that the accurate validation of results is paramount for successful missions. This dedication ensures that when PythonProgramming_Expert presents a solution, it is not only theoretically sound but also practically reliable."

11 },

We also summarized the agent-selected rate for each scenario in Figure 5. It is obvious that selected agents are highly related to the current scenarios. The verification expert has a high selection rate because we prompt Captain Agent in the system message to create a verification role to maintain the conversation. We also notice that in some specific scenarios (mathematics, data analysis, and programming), some agents with a general name and description will have a high selection rate (e.g., PythonMath_Expert, DataAnalysis_Expert, CodeReview_Expert, etc.). However, in the Science scenarios, there are no highly preferred agents with general descriptions, and the selection distribution become flatten.

1242 H TOOL LIBRARY

1243
1244 This section provides the names and descriptions of our manually created tool library. The tools are
1245 categorized into three classes: Information Retrieval, Data Analysis and Math Problem Solving. For
1246 each category, we summarize the patterns of the corresponding dataset and manually craft a set of
1247 functions suits the tasks and can potentially enhance the agents' task resolution capability.
1248

1249
1250 Table 8: Tools for Information Retrieval category.

1251 Tools	1251 Description
1252 scrape_wikipedia_tables	1252 Scrapes Wikipedia tables based on a given URL and header key- 1253 word.
1254 transcribe_audio_file	1254 Transcribes the audio file located at the given file path.
1255 youtube_download	1255 Downloads a YouTube video and returns the download link.
1256 academic_search	1256 Perform an academic search of papers, authors or an author's 1257 papers.
1258 docx_to_md	1258 Converts a DOCX file to Markdown format.
1259 pptx_to_md	1259 Convert a PowerPoint presentation (PPTX) to Markdown format.
1260 spreadsheet_to_md	1260 Convert an Excel spreadsheet file to Markdown format.
1261 extract_pdf_image	1261 Extracts images from a PDF file and saves them to the specified 1262 output directory.
1263 extract_pdf_text	1263 Extracts text from a specified page or the entire PDF file.
1264 get_youtube_caption	1264 Retrieves the captions for a YouTube video.
1265 image_qa	1265 Answers your questions about a given image.
1266 optical_character_recognition	1266 Perform optical character recognition (OCR) on the given image.
1267 perform_web_question_answering	1267 Perform web search according to keyword and answer your ques- 1268 tion on each webpage search result, or directly on the webpage if 1269 the keyword is a URL. For each search result, a response to the 1270 question is provided.
1271 scrape_wikipedia_tables	1271 Scrapes Wikipedia tables based on a given URL and header key- 1272 word. 1273

1274
1275
1276 Table 9: Tools for Data Analysis category.

1277 Tools	1277 Description
1278 calculate_correlation	1278 Calculate the correlation between two columns in a CSV file.
1279 calculate_skewness_and_kurtosis	1279 Calculate the skewness and kurtosis of a specified column in a 1280 CSV file. The kurtosis is calculated using the Fisher definition.
1281 detect_outlier_iqr	1281 Detect outliers in a specified column of a CSV file using the IQR 1282 method.
1283 detect_outlier_zscore	1283 Detect outliers in a CSV file based on a specified column. The 1284 outliers are determined by calculating the z-score of the data 1285 points in the column.
1286 explore_csv	1286 Reads a CSV file and prints the column names, shape, data types, 1287 and the first few lines of data.
1288 shapiro_wilk_test	1288 Perform the Shapiro-Wilk test on a specified column of a CSV 1289 file.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Table 10: Tools for Math Problem solving category.

Tools	Description
calculate_circle_area_from_diameter	Calculate the area of a circle given its diameter.
calculate_day_of_the_week	Calculates the day of the week after a given number of days starting from a specified day.
calculate_fraction_sum	Calculates the sum of two fractions and returns the result as a mixed number.
calculate_matrix_power	Calculate the power of a given matrix.
calculate_reflected_point	Calculates the reflection point of a given point about the line $y=x$.
complex_numbers_product	Calculates the product of a list of complex numbers.
compute_currency_conversion	Compute the currency conversion of the given amount using the provided exchange rate.
count_distinct_permutations	Counts the number of distinct permutations of a sequence where items may be indistinguishable.
evaluate_expression	Evaluates a mathematical expression with support for floor function notation and power notation.
find_continuity_point	Find the value that ensures the continuity of a piecewise function at a given point.
fraction_to_mixed_numbers	Simplifies a fraction to its lowest terms and returns it as a mixed number.
modular_inverse_sum	Calculates the sum of modular inverses of the given expressions modulo the specified modulus.
simplify_mixed_numbers	Simplifies the sum of two mixed numbers and returns the result as a string in the format 'a b/c'.
sum_of_digit_factorials	Calculates the sum of the factorial of each digit in a number.
sum_of_primes_below	Calculates the sum of all prime numbers below a given threshold.