MAS-GPT: TRAINING LLMS TO BUILD LLM-BASED MULTI-AGENT SYSTEMS

Anonymous authors

Paper under double-blind review

ABSTRACT

LLM-based multi-agent systems (MAS) have shown significant potential in tackling diverse tasks. However, to design effective MAS, existing approaches heavily rely on manual configurations or multiple calls of advanced LLMs, resulting in inadaptability and high inference costs. In this paper, we simplify the process of building an MAS by reframing it as a generative language task, where the input is a user query and the output is a corresponding MAS. To address this novel task, we unify the representation of MAS as executable code and propose a consistencyoriented data construction pipeline to create a high-quality dataset comprising coherent and consistent query-MAS pairs. Using this dataset, we train MAS-GPT, an open-source medium-sized LLM that is capable of generating query-adaptive MAS within a single LLM inference. The generated MAS can be seamlessly applied to process user queries and deliver high-quality responses. Extensive experiments on 9 benchmarks and 4 LLMs show that the proposed MAS-GPT consistently outperforms 10+ baseline MAS methods on diverse settings, indicating MAS-GPT's high effectiveness, efficiency and strong generalization ability.

025 026

027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Large language models (LLMs) such as ChatGPT (Ouyang et al., 2022; OpenAI, 2023) have
achieved significant success on a wide range of tasks. However, a single LLM often struggles to
handle the diverse and complex range of tasks (e.g., varying difficulties and domains) encountered
in practice (Hong et al., 2024; Chen et al., 2024).

Such limitation has driven recent research towards building LLM-based multi-agent systems (MAS) (Qian et al., 2024a; Chen et al., 2024; Liu et al., 2024b), where multiple LLMs (agents) with specialized capabilities work collaboratively to achieve more effective solutions. For example, MetaGPT (Hong et al., 2024) and ChatDev (Qian et al., 2024a) build multi-LLM teams with expertise roles (e.g., programmer, tester, and product manager) to solve complex coding tasks in a predefined pipeline; while AgentVerse (Chen et al., 2024) involves recruiters, executors, and evaluators for iterative task solving. These methods have shown superior performance over single LLM inference.

040 Despite achieving promising task performance, there are two fundamental issues that hinder the 041 broad applications of MAS: inadaptability and high costs. (i) Inadaptability & high human effort: 042 MAS in MetaGPT (Hong et al., 2024), ChatDev (Qian et al., 2024a), and AgentVerse (Chen et al., 043 2024) are all manually crafted (e.g., for coding tasks). That is, the collaboration structure and 044 agents' prompts are predetermined and static, lacking in the generality to adapt towards any given 045 tasks. (ii) High inference costs: Although there have been efforts to design adaptive MAS, they essentially shift the human cost onto the computational cost. For example, both GPTSwarm (Zhuge 046 et al., 2024) and DyLAN (Liu et al., 2024b) rely on LLMs to replace human involvement, iteratively 047 adjusting the collaboration structure or agents' prompts in the MAS for each specific task. However, 048 this process often requires multiple LLM inferences. 049

Focusing on these key issues, this paper explores how to adaptively build a query-specific MAS at
a minimal cost. Our core idea is to reframe the process of building an executable MAS for each
query as a generative language task, making building MAS as simple and efficient as querying ChatGPT (Ouyang et al., 2022). Given the generated MAS, the query can then be seamlessly processed
to produce the final response, significantly simplifying the whole pipeline.

Under this context, we introduce MAS-GPT, 055 an LLM specifically trained to adaptively gen-056 erate executable MAS based on any given user 057 query in one single inference. While the con-058 cept is straightforward, the challenge lies in the limited knowledge of LLMs on the task of MAS generation and the lack of corresponding 060 training data. These limitations raise two key 061 technical challenges: how to represent the MAS 062



Figure 1: Introduction of our proposed new paradigm for building MAS. During inference, MAS-GPT adaptively generates a query-specific MAS with one LLM inference.

and how to construct the dataset. (1) To ensure the generated MAS is readily executable, we unify 063 the representation of MAS by describing it as a Python code snippet (i.e., a forward function), with 064 each agent's prompt as a variable, LLM calls as functions, and agent relationships as string con-065 catenation. (2) Building on this foundation, we propose a consistency-oriented data construction 066 pipeline to facilitate the model in learning generalizable patterns and logical correlations, which in-067 cludes the construction, evaluation, selection, and refinement of query-MAS pairs. During selection, 068 we design an inter-consistency-oriented selection approach to ensure that similar queries are paired with similar high-performing MAS, facilitating the model to learn generalizable patterns. During 069 refinement, we propose a intra-consistency-oriented refinement method to strengthen the relatedness between query and MAS, enabling the model to learn the reasonable correlation. Finally, the result-071 ing pairs are used to train open-source LLMs via supervised fine-tuning, where the instruction is the 072 user query and the response is the MAS represented by code. This will equip the model with the 073 ability to generate query-specific MAS, and also, generalize to unseen queries. 074

With the introduction of MAS-GPT, inference for a query becomes significantly simplified. Instead of relying on manual crafting (Hong et al., 2024; Qian et al., 2024a; Chen et al., 2024) or multiple LLM inference costs (Liu et al., 2024b; Zhuge et al., 2024) to obtain an MAS for each query, the user simply inputs a query into MAS-GPT to get a corresponding executable MAS. Such MAS can be directly applied to process the query, where multiple MAS-GPT-generated agents collaborate with an MAS-GPT-generated structure to deliver the final solution. With advantages of adaptability, low cost, and generalization, this approach could facilitate the application of MAS at scale.

We conduct extensive experiments to compare MAS-GPT with 10+ baseline methods on 9 benchmarks (various domains) using 4 state-of-the-art (open-source and proprietary) LLMs. Our results show that MAS-GPT consistently outperforms baseline methods on average, indicating its high generality and effectiveness. Meanwhile, MAS-GPT has the potential to further push the boundary of strong reasoning capability of o1 (OpenAI, 2024b), bringing 13.34% gain on AIME-2024, a challenging mathematical benchmark. We also verify that our MAS-GPT can generalize to unfamiliar queries and generate novel MAS via case studies.

- 089 Our contributions are as follows:
- We reframe building MAS for each query as a generative language task. We unify the representation of MAS as executable code and propose a consistency-oriented query-MAS data construction pipeline for LLM training.
 - 2. We introduce MAS-GPT, an LLM specifically trained to generate query-specific executable MAS. All code, data and models will be open-sourced.
 - 3. Experiments on 9 benchmarks and 4 LLMs show that MAS-GPT consistently outperforms 10+ baselines, indicating its effectiveness, efficiency and generalization ability.
- 098 099

100

094

095

096

2 RELATED WORK

LLM-based Multi-Agent Systems. Since a single LLM may struggle to handle the diverse and complex range of tasks in practice (Li et al., 2023; Qian et al., 2024b), such limitation has driven recent research towards building LLM-based multi-agent systems (MAS) (Wang et al., 2024c; Wu et al., 2023). MetaGPT (Hong et al., 2024) and ChatDev (Qian et al., 2024a) introduce manually designed multi-agent teams for solving coding tasks; while MedAgents is designed for medical tasks (Tang et al., 2024). AgentVerse (Chen et al., 2024) proposes an iterative collaboration structure where agents are recruited to discuss, execute, and evaluate. Multi-Agent Debate (Du et al., 2024; Liang et al., 2024) designs multiple expertise LLM-agents to debate and reason over multiple rounds

to get final answers. The MAS in these methods are all fixed regardless of the given query, lacking
 in the generality to adapt accordingly.

110 DyLAN (Liu et al., 2024b) leverages LLMs to evaluate agents' values and dynamically select the 111 best agents, GPTSwarm (Zhuge et al., 2024) manually initializes an agent team, adjusts the collabo-112 ration structure and agents' prompts by prompting LLMs. Given queries with ground-truth answers 113 from one task and several available MAS as context, ADAS (Hu et al., 2024) and AFlow (Zhang 114 et al., 2024) leverages the strong capabilities of LLMs such as Claude-3.5-sonnet (Anthropic, 2024) 115 and GPT-4 (OpenAI, 2023) to iteratively generate task-oriented MAS for the specific task. All these 116 methods require multiple times of LLM calls (e.g., over 10 calls of API with lengthy context) to ob-117 tain an MAS for each specific query, which is time-consuming and compute-expensive in practice.

Instead of manually designing a fixed MAS (Qian et al., 2024a; Chen et al., 2024; Du et al., 2024) or requiring multiple LLM inference costs to obtain an MAS (Liu et al., 2024b; Zhuge et al., 2024) for each query, our MAS-GPT significantly simplifies the process of building an MAS, which can flexibly generate query-specific MAS within one LLM inference. Specifically, we design a data-construction pipeline to generate a series of query-MAS pairs, which are used for training MAS-GPT based on open-source LLMs.

124 **LLM Post-Training.** Modern state-of-the-art LLMs are usually post-trained via two main stages: 125 supervised fine-tuning (SFT) and preference learning (Ouyang et al., 2022; Dubey et al., 2024; Yang 126 et al., 2024; Liu et al., 2024a), where SFT is the basic technique to teach LLM a defined tasks (Zhou 127 et al., 2023; Longpre et al., 2023). Focusing on SFT, a series of researches are conducted on the 128 construction of datasets for training chatbot-type LLMs. For example, LIMA (Zhou et al., 2023) 129 manually annotates high-quality language data for SFT, emphasizing the importance of quality of 130 SFT datasets. WizardLM (Xu et al., 2024), TULU 3 (Lambert et al., 2024), and Persona Hub (Ge 131 et al., 2024) synthesize SFT data by prompting GPT models, indicating the potential of synthetic data for LLM training. For MAS-GPT, the training process leverages SFT, with a primary focus on data 132 construction. While previous approaches focus on training LLMs to directly answer user queries, 133 the challenge of training LLMs to generate MAS from user queries introduces a novel difficulty. 134 Unlike real-world dialogue data, LLMs have limited (if any) knowledge of MAS generation. Using 135 our proposed data construction pipeline, we create the first query-MAS-paired dataset, which will 136 be made open-source in future. 137

137 138 139

140

3 Methodology

This section first outlines the overall system integrated with MAS-GPT when processing user queries during inference. Next, we delve into the specifics of training MAS-GPT, with a particular focus on the dataset construction process.

144 145

3.1 OVERALL SYSTEM INTEGRATED WITH MAS-GPT

We follow a standard workflow: given a user query, a multi-agent system (MAS) is constructed, with
multiple agents working collaboratively to generate the final answer. Unlike previous approaches
that either manually design the MAS, rely on fixed and query-agnostic MAS, or incur significant
computational costs to determine the appropriate MAS, our approach streamlines the entire process
of building MAS by reducing it to a single LLM inference.

The core of our system is MAS-GPT, an LLM that is trained to generate MAS tailored specifically to the input query. Instead of relying on pre-built agent configurations, MAS-GPT dynamically creates an MAS for each query, ensuring that the system adapts to a wide range of tasks. This approach not only minimizes the time and computational resources traditionally required to build the right MAS but also enhances the system's flexibility by generating task-specific solutions in real-time. Finally, the MAS generated by MAS-GPT can be seamlessly integrated to process the query and deliver the final answer (bottom right in Figure 3).

158 159

160

3.2 MAS-GPT: DATASET CONSTRUCTION AND TRAINING

161 To achieve the above goal, we reframe building MAS as a generative language task, where the input is a user query and the output is an executable MAS capable of processing that query. This shift to a



Figure 3: Illustrations of dataset construction, training, and inference of our proposed MAS-GPT.

generative paradigm introduces a new challenge since there is few (if any) knowledge within LLMs
on MAS generation. To make this approach viable, the key focus lies in constructing an appropriate
dataset to teach the LLMs such brand-new task. To achieve this, we propose a consistency-oriented
data construction process, which involves four key steps: (1) construction of query and MAS pools,
(2) inference and evaluation of query-MAS pairs, (3) inter-consistency-oriented pair selection, (4)
intra-consistency-oriented pair refinement.

183 Data - Construction of Query and MAS Pools (Representing MAS as Executable Code). To construct the dataset for supervised fine-tuning (SFT), we adopt the following data format: (system 184 prompt, instruction, response). Here, the system prompt briefly describes the MAS generation task, 185 the instruction corresponds to the user query, and the response includes the MAS, which can be extracted by string matching. Therefore, training the LLM requires the collection of a series of 187 query-MAS pairs. Firstly, to enable MAS-GPT to handle diverse queries, we build a query pool 188 from open-source queries across various domains, such as general QA, mathematics, and coding. 189 Each query is carefully selected to be verifiable, ensuring the presence of a ground-truth answer or 190 test cases (e.g., for coding tasks). 191

While the collection of queries is relatively 192 straightforward, constructing the MAS pool 193 presents a fundamental challenge: how to rep-194 resent an executable MAS. To address this, we 195 propose unifying the representation of MAS by 196 formalizing it as executable Python code snip-197 pets. This unified representation is motivated 198 by the observation that all existing LLM-based 199 MAS methods are ultimately implemented as 200 code, encompassing the definition of agents' 201 prompts, LLM calls, and inter-agent interactions (Qian et al., 2024a; Hu et al., 2024; Zhang 202 et al., 2024). Specifically, we define an MAS as 203 a forward function that takes a user query as in-204

174 175 176



Figure 2: Our unified code representation of an executable MAS (i.e., a forward function). Each color denotes an agent. Agents defined by variables, LLM calls denoted by function calls, and interactions represented by string concatenations.

put and returns the final answer generated by the MAS. Within the forward function, agent prompts are defined as variables, agent inferences are implemented as function calls, and interactions between agents are represented through string concatenation; see an example in Figure 2.

Following this framework, we first re-implement several existing MAS methods (e.g., Multi-Agent Debate (Du et al., 2024), Self-Consistency (Wang et al., 2024b), Self-Refine (Madaan et al., 2024))
to align with our unified code representation. To further expand the diversity of MAS candidates, we also manually design some MAS systems, resulting in a base MAS pool comprising over 40 unique MAS designs (Figure 7). Importantly, these 40+ MAS do not directly correspond to the exact number of MAS in the training dataset; rather, they serve as foundations that evolve during the query-MAS pair refinement process.

Data - Evaluation of Query-MAS Pairs. After constructing the query and MAS pools, it is crucial to evaluate the query-MAS compatibility since not all MAS designs are equally suitable for every

query. To achieve this, we pair each query and MAS in the pool by inferring the query to the MAS and evaluating the generated final answer.

Specifically, given the query pool $\mathbf{Q} = \{(Q_i, Y_i)\}_{i=1}^N$ and the base MAS pool $\mathbf{M} = \{MAS_i^{\text{base}}\}_{i=1}^M$. 219 where Q_i is the query, Y_i is the information for verification, N and M denotes the pool size, 220 we obtain $N \times M$ pairs. Then, a query-dependent evaluation function $f_{\text{eval}}(\cdot)$ will be applied 221 to evaluate the effectiveness of the query-MAS pair: $score_{i,j} = f_{eval}(MAS_j^{base}(Q_i), Y_i)$, where 222 $MAS_{i}^{\text{base}}(Q_{i})$ denotes the answer generated by MAS_{i}^{base} given the query Q_{i} , 1 and 0 denotes cor-223 rect and wrong respectively. Overall, we get M MAS scores for each query Q_i , which are denoted 224 by $\mathbf{s}_i = [score_{i,1}, ..., score_{i,M}]$, laying the foundation for subsequent steps for selecting appropriate 225 query-MAS pairs and further refinement. 226

Data - Inter-Consistency-Oriented Pair Selection. With the query-MAS pair results obtained
 from the evaluation step, the next critical task is to select and construct high-quality query-MAS
 pairs for training. The first selection criterion is intuitive: *effectiveness*. Specifically, we retain
 only the query-MAS pairs where the MAS produces a correct answer (evaluation score is 1), as
 MAS designs that generate correct answers are more likely to be suitable for their respective queries
 compared to those that fail.

While using all the remaining effective query-MAS pairs for training is straightforward, it introduces a significant problem of low *inter-consistency*: the same or similar queries may correspond to multiple different MAS designs. This lack of consistency makes it difficult for the model to learn a clear optimization objective, hindering its ability to understand and perform the task effectively.

To address this issue, we propose an inter-consistency-oriented pair selection method that optimizes both *effectiveness* and *inter-consistency*. The core idea is to group similar queries and assign them a single, high-performing MAS to maintain consistency across the dataset. Specifically, we cluster queries based on their metadata or embeddings. For a group of S queries $\mathbb{S} = \{Q_i\}_{i=1}^S$, we calculate a cumulative score for each MAS by summing its effectiveness scores across all queries in the group: $\mathbf{s} = \sum_{i=1}^{S} \mathbf{s}_i$. The MAS with the highest cumulative score is then selected as the representative MAS for all queries in the group: $MAS_{ass}^{base} = \arg \max_{MAS \in \mathbf{M}} \mathbf{s}$. Through this, we pair each specific query with a specific base MAS: (Q_i, MAS_i^{base}) .

By aligning similar queries with the same high-performing MAS, it improves the inter-consistency
of the query-MAS pairs, helping the model recognize generalizable patterns and generalize across
similar queries. For example, queries requiring divergent thinking may be consistently paired with
MAS structures where multiple agents independently generate ideas and then discuss.

Data - Intra-Consistency-Oriented Pair Refinement. While the inter-consistency-oriented pair
 selection process effectively ensures consistency across query-MAS pairs, there remains a critical
 issue within individual pairs: *intra-consistency*. Specifically, the alignment between a query and its
 associated MAS may still be suboptimal, making it challenging for the model to learn meaningful
 associations. For instance, a query about physics may be paired with an MAS involving experts
 from multiple domains (e.g., physics and biology), where the presence of non-relevant agents like
 biology experts can confuse the model.

256 To address this, we propose an intra-consistency-oriented pair refinement method. This approach 257 aims to improve the query-MAS alignment through two key strategies: (1) adjusting MAS to make it query-dependent, and (2) introducing a reasoning process to strengthen the connection between 258 the query and MAS. We employ an LLM-based data synthesis method, where an advanced LLM 259 adjusts agents' definitions within the MAS based on the query and the previously selected MAS. 260 The LLM is also instructed to generate a reasoning statement that explains the relationship between 261 the query and the refined MAS, improving the interpretability of the query-MAS pair; please refer 262 to prompt in Table 10. This process enables the model to better understand the context and logic 263 behind each decision, which in turn facilitates model training and improves generalization. 264

Next, we infer and evaluate the refined MAS on the corresponding query, as advanced LLMs could generate inappropriate or non-executable MAS. Specifically, for each base pair $(Q_i, MAS_i^{\text{base}})$, the refined MAS_i^{refine} is tested and accepted only if it achieves a not-worse score. Formally:

$$MAS_{i} = \begin{cases} MAS_{i}^{\text{refine}}, & \text{if } s^{\text{refine}} \geq s^{\text{base}} \\ MAS_{i}^{\text{base}}, & \text{otherwise} \end{cases},$$

Table 1: Comparing MAS-GPT with 10 baselines across 8 benchmarks using Llama-3-70B-Instruct, MAS-GPT performs the best on average, verifying its generality in handling diverse queries. Benchmarks with * are out-of-domain for MAS-GPT.

Method	MATH	GSM8K	GSM-H	H-Eval*	H-Eval+*	MMLU	GPQA*	SciBench*	Avg.
Single (Dubey et al., 2024)	50.55	92.38	45.80	79.01	75.78	77.37	36.68	21.05	59.83
Chain-of-Thought (Wei et al., 2022)	53.20	92.79	46.20	77.16	77.02	75.56	35.28	17.68	59.36
Self-Consistency (Wang et al., 2024b)	61.59	94.99	47.20	77.78	75.78	78.18	37.15	20.00	61.58
LLM-Debate (Du et al., 2024)	61.37	91.58	44.60	74.69	74.53	77.78	34.35	19.79	59.84
Self-Refine (Madaan et al., 2024)	58.50	90.78	37.80	67.90	62.73	74.75	38.32	20.00	56.35
Quality-Diversity (Lu et al., 2024)	60.49	92.99	45.60	70.99	70.19	75.76	33.64	20.63	58.79
SPP (Wang et al., 2024c)	51.66	92.79	44.80	76.54	73.29	77.37	35.05	20.84	59.04
AgentVerse (Chen et al., 2024)	55.63	93.39	41.40	77.78	73.91	76.57	40.19	16.00	59.36
GPTSwarm (Zhuge et al., 2024)	55.41	93.19	43.20	69.14	73.91	75.15	36.45	14.11	57.57
DyLAN (Liu et al., 2024b)	59.60	91.18	44.80	79.01	75.78	78.18	35.98	19.79	60.54
MAS-GPT (Ours)	68.65	93.39	62.40	80.25	78.88	78.38	37.62	24.21	65.47

283

where $s^{\text{refine}} = f_{\text{eval}}(MAS_i^{\text{refine}}(Q_i), Y_i)$ and $s^{\text{base}} = f_{\text{eval}}(MAS_i^{\text{base}}(Q_i), Y_i)$ are evaluation scores by comparing the MAS-generated and ground-truth answer Y_i .

Through this process, each query Q_i is ultimately associated with a tuple (Q_i, R_i, MAS_i) , where R_i denotes the reasoning statement, and MAS_i is the final MAS. This refined dataset ensures both inter- and intra-consistency, providing high-quality training data for subsequent model fine-tuning.

Training - Supervised Fine-Tuning of MAS-GPT Our dataset follows the format (system prompt, instruction, response). The system prompt briefly describes the task of generating a query-specific MAS and the instruction corresponds to the user query Q_i . The response is constructed as the concatenation of the reasoning process and the final MAS, which is represented as executable code in text form.

Building upon this dataset, we perform supervised fine-tuning of MAS-GPT on the open-source medium-sized LLM, Qwen2.5-Coder-32B-Instruct (Yang et al., 2024), leveraging its capabilities of code generation and instruction-following. During inference, when a user query is received, MAS-GPT generates an executable MAS tailored to that specific query Q_i : $MAS_i^{\text{gen}} = \text{MAS-GPT}(Q_i)$. The generated MAS is directly usable for processing the query Q_i and delivering the final answer: $A_i = MAS_i^{\text{gen}}(Q_i)$, significantly simplifying the task handling process.

3.3 DISCUSSIONS

300 301

299

Overall, our system integrated with MAS-GPT offers the following key advantages: simplicity, cost-302 efficiency, and adaptability (generality). Instead of manually designing an MAS for each specific 303 query, relying on a fixed MAS for all queries, or requiring multiple LLM inference costs to obtain an 304 MAS for a query, our MAS-GPT significantly simplifies the process of building an MAS by reducing 305 into one single LLM inference. Given a user query, MAS-GPT will efficiently return a query-specific 306 MAS, which is executable and can be seamlessly applied to process the query to deliver the final answer. Although training incurs some cost, it is a one-time expense, whereas inference is potentially 307 endless in practical applications. We believe that MAS-GPT has the potential to further advance the 308 real-world application of MAS due to its simplicity, cost-efficiency, and adaptability. 309

- 310
- 311 4 EXPERIMENTS
- 312 313

314

4.1 EXPERIMENTAL SETUPS

Training. Our training queries are sampled from the training splits available in MATH (Hendrycks 315 et al., 2021b), GSM8K (Cobbe et al., 2021), MBPP (Austin et al., 2021), MMLU (Hendrycks et al., 316 2021a), and SciQ (Welbl et al., 2017), covering domains of math, coding, and general QA. Llama-317 3-70B-Instruct is used during dataset construction. The number of training samples (i.e., query-318 MAS pairs) is approximately 11k; see the statistics of our dataset in Table 4. Our MAS-GPT is 319 trained over Qwen2.5-Coder-32B-Instruct (Yang et al., 2024), leveraging its instruction-following 320 and coding capabilities. We train the LLM using 16 A100s with an effective batch size of 32 for 3 321 epochs at a learning rate of 1e-5 (Zheng et al., 2024). 322

Testing. To verify that our MAS-GPT can handle diverse queries in practice, we consider multiple benchmarks from diverse domains. These include MATH (Hendrycks et al.,



Figure 4: (a) Different methods empowered with strong reasoning LLM: o1-preview. We see that our MAS-GPT significantly enhance the reasoning performance over single LLM, indicating its potential in further augmenting LLM reasoning. (b) Comparisons with AFlow (optimized on MATH).
MAS-GPT even outperforms AFlow on its in-domain benchmarks; while AFlow fails on out-of-domain benchmarks. (c) MAS-GPT achieves the best performance with low inference cost.

2021b), GSM8K (Cobbe et al., 2021), and GSM-Hard (Gao et al., 2023) for math domains; HumanEval (Chen et al., 2021) and HumanEval+ (Liu et al., 2023) for coding tasks;
MMLU (Hendrycks et al., 2021a) for general QA tasks; GPQA (Rein et al., 2023) and
SciBench (Wang et al., 2024a) for science topics. Please refer to Table 9 for details about datasets and Section C.2 for details about evaluation. For all baselines, the LLMs that drive the MAS to process user queries are kept the same, where we consider four state-of-the-art LLMs including Llama-3-70B-Instruct (Dubey et al., 2024), Qwen2.5-72B-Instruct (Yang et al., 2024), GPT-4o-mini-2024-07-18 (OpenAI, 2024a), and o1-preview-2024-09-12 (OpenAI, 2024b).

345 **Baselines.** For fair comparisons, we consider 10 baselines that are suitable for handling diverse 346 tasks. We include single agent and agent with chain-of-thought (Wei et al., 2022) as two ba-347 sic baselines, Self-Consistency (Wang et al., 2024b) and Quality-Diversity (Lu et al., 2024) that 348 select the best from multiple answers, LLM-Debate (Du et al., 2024) that involves multiple ex-349 perts for debating, Self-Refine (Madaan et al., 2024) that iteratively refines last agent's answer, 350 SPP (Wang et al., 2024c) that stimulates conversations among multiple roles, AgentVerse (Chen 351 et al., 2024) and DyLAN (Liu et al., 2024b) that dynamically adjust multi-agent team during infer-352 ence, GPTSwarm (Zhuge et al., 2024) that relies on a graph collaboration structure.

353 354

355

333

334

335

336

337

4.2 MAIN RESULTS

Since our MAS-GPT aims to facilitate the multi-agent systems in flexibly handling diverse queries,
 our results focus on generality. Here, we show the generality of MAS-GPT by comparing performance averaged on various benchmarks and performance using different LLMs to drive the MAS.

MAS-GPT's generality in handling diverse queries. We compare MAS-GPT with 10 baselines
 on 8 benchmarks using Llama-3-70B-Instruct (Dubey et al., 2024) to drive the MAS, with results
 reported in Table 1. GPQA and SciBench are two benchmarks that are out-of-domain for our MAS GPT. From the table, we see that (1) our MAS-GPT significantly outperforms the baseline methods
 on average, outperforming the second-best method by 3.89%. (2) Our MAS-GPT simultaneously
 achieves promising performance in both in-domain and out-of-domain (i.e., queries that are significantly different from those in the training data) benchmarks, indicating MAS-GPT's generality.

366 Generality in using diverse LLM backbones for MAS. Llama-3-70B-Instruct was utilized to 367 drive MAS during dataset construction for training MAS-GPT, a 32B-sized LLM. As shown in 368 Table 1, this approach proves effective when employing the same LLM to drive MAS during test time. To further validate the versatility of MAS-GPT, we assess its performance under different 369 MAS-driving LLMs, including Qwen2.5-72B-Instruct and GPT-4o-mini-2024-07-18, in Table 2. 370 The results demonstrate that MAS-GPT consistently achieves superior performance, regardless of 371 the LLM used to drive MAS, highlighting its strong compatibility and adaptability across various 372 MAS-driving LLMs. 373

MAS-GPT's potential in further augmenting the reasoning performance of strong reasoning
 LLMs such as o1. In recent developments, the AI community has introduced several state-of-the art reasoning LLMs (OpenAI, 2024b; Qwen, 2024), which have demonstrated remarkable reason ing capabilities by scaling inference-time computations (Snell et al., 2024). In this context, we
 aim to explore whether our proposed MAS-GPT can take the reasoning power of these already ad-



Figure 5: Explorations of scaling in training MAS-GPT. (a) More data leads to fewer execution failures. (b) More data contributes to better performance of MAS-GPT in facilitating MAS application. Without training (N=0), the model fails, highlighting that MAS generation is a non-trivial task requiring specific training. (c) Larger model generally contributes to better performance. These findings demonstrate the promising potential of MAS-GPT, suggesting that it can be further improved with more diverse, high-quality data and stronger models as the community continues to advance.

vanced models even further. To test this, we conduct experiments using OpenAI's o1-preview-202409-12 (OpenAI, 2024b) model, evaluating it on the highly challenging AIME-2024 mathematical
benchmark¹. The results, as shown in Figure 4(a), show that our proposed MAS-GPT significantly
outperforms the baseline methods on this challenging task. Specifically, it improves over the single LLM by a large margin: 13.34%. This result not only verifies the generality of our proposed
MAS-GPT, but also indicates its promising potential in pushing the boundaries of LLM reasoning.

Comparisons with task-specific methods, 399 **AFlow.** To further demonstrate the generality 400 and effectiveness of our MAS-GPT during in-401 ference time, we compare with AFlow (Zhang 402 et al., 2024), a latest task-specific method for 403 MAS optimization that has been specifically 404 optimized on MATH (Hendrycks et al., 2021b) 405 We evaluate on two AFlow's indataset. 406 domain (MATH and GSM8K) and two AFlow's out-domain (MMLU and HumanEval+) bench-407 marks. Results in Figure 4(b) show surpris-408 ingly good performance of our proposed MAS-409 GPT. As a general method, our MAS-GPT even 410 outperforms math-specific AFlow on the MATH 411 dataset by 3.53%! Meanwhile, the MAS opti-412 mized on MATH by AFlow fails to generalize 413 to other domains, achieving worse performance 414 than a single LLM. In contrast, our MAS-415 GPT consistently performs the best across these 416 benchmarks. It is also worth mentioning that 417 our MAS-GPT only requires one-time infer-418 ence of a 32B-sized LLM to build the MAS; while AFlow needs to call the APIs of pow-419 erful proprietary LLMs, such as Claude-3.5-420

Table 2: MAS-GPT consistently performs the best across MAS-driving LLMs, indicating its strong compatibility.

Method	MATH	GSM-H	H-Eval+	MMLU	GPQA	Avg.
	C)wen2.5-7	72B-Instru	uct		
Single	85.86	64.91	85.37	82.60	44.39	72.63
COT	86.90	62.27	84.15	83.20	47.86	72.88
Self-Con.	87.32	61.46	87.20	83.40	50.00	73.88
LLM-Debate	85.24	63.49	68.90	86.20	47.86	70.34
Self-Refine	83.58	59.03	78.66	85.40	43.32	70.00
Q-D	85.65	63.08	76.83	82.80	48.66	71.40
SPP	85.65	62.88	82.32	83.40	48.40	72.53
AgentVerse	84.82	59.43	81.10	83.20	44.65	70.64
GPTSwarm	83.16	63.89	83.54	84.60	44.92	72.02
DyLAN	87.73	63.08	85.37	84.40	51.07	74.33
MAS-GPT	87.53	66.33	85.98	83.80	48.66	74.46
	G	PT-4o-mi	ni-2024-0′	7-18		
Single	78.18	58.03	86.25	78.56	38.03	67.81
COT	78.79	60.84	85.62	79.16	39.60	68.80
Self-Con.	81.62	59.04	85.00	80.96	39.82	69.29
LLM-Debate	79.60	60.84	86.25	80.76	37.81	69.05
Self-Refine	74.55	54.62	76.88	79.16	33.33	63.71
Q-D	79.80	59.64	84.38	79.76	37.58	68.23
SPP	77.58	57.63	86.25	77.96	37.58	67.40
AgentVerse	75.15	55.62	79.38	78.36	36.24	64.95
GPTSwarm	75.15	55.62	79.38	78.36	36.32	64.97
DyLAN	81.21	59.24	80.62	79.96	40.94	68.39
MAS-GPT	81.21	61.45	86.88	80.36	42.60	70.50

Sonnet (Anthropic, 2024), 10 times per query and depends on a hold-out validation set.

Cost comparisons. Here, we compare the inference cost of various methods from the moment a
 user query is received to the generation of the final answer, as illustrated in Figure 4(c). We quantify
 the inference cost in terms of the number of LLM inference calls (Liu et al., 2024b), interpreting the
 inference of MAS-GPT as 0.5 times, given that its model size is approximately half that of the MAS driving LLM (32B v.s. 70B). From the figure, we observe that, among the four methods compared,
 MAS-GPT achieves the best performance while requiring the fewest inference calls, demonstrating
 its efficiency and effectiveness.

429 430

431

387

389

390

391

392

4.3 ANALYSIS OF MAS-GPT

¹https://huggingface.co/datasets/Maxwell-Jia/AIME_2024

432 Effectiveness of inter-consistency-oriented 433 pair selection. During data construction, to 434 facilitate the model in recognizing generaliz-435 able patterns between queries and MAS, we propose an inter-consistency-oriented query-436 MAS pair selection method, which maps sim-437 ilar queries with consistent high-performing 438 MAS. To examine its effectiveness, we re-439 place this mapping with a random mapping 440 approach, which randomly selects one out of 441 those MAS with correct answers. From Ta-442 ble 3, by comparing 1 and 4, we see that 443 our proposed method brings significant per-444 formance gain, with an absolute improvement 445 of 8.39% on MATH.

Table 3: Ablation studies on the designs of dataset construction: (1) our inter-consistency-oriented pair selection, (2) the adjustment of MAS in our intraconsistency-oriented pair refinement: Refine-A, (3) the introduction of reasoning process in our intraconsistency-oriented pair refinement: Refine-R. The table shows that these three designs all play critical roles in achieving high task performance.

	Select	Refine-A	Refine-R	MATH	MMLU	GPQA
1	×	\checkmark	\checkmark	60.26	77.58	36.68
2	\checkmark	×	\checkmark	66.23	77.78	36.45
3	\checkmark	\checkmark	X	64.90	75.96	37.15
4	\checkmark	\checkmark	\checkmark	68.65	78.38	37.62

446 Effectiveness of intra-consistency-oriented pair refinement. During data construction, to help 447 the model learn the associations between query and MAS, we propose an intra-consistency-oriented 448 query-MAS pair refinement method. This method enhances the query-MAS alignment by adjusting 449 MAS to make it query-dependent and introducing a reasoning process to strengthen the logical 450 connection. To examine their effects, we conduct two experiments with one without adjustment of 451 MAS and one without reasoning process. From Table 3, by comparing 2 and 4, 3 and 4, we see 452 that our designs in adjusting MAS and introducing reasoning process both contribute to performance 453 improvement, indicating the effectiveness of our proposed refinement method.

454 Scaling effects of data size. To explore the scaling effects of data size for training MAS-GPT, we 455 adjust the size from 0 to 11k using the same 32B-sized model and compare the extractability (i.e., the 456 Python code can be extracted), executability (i.e., the code is executable), task performance. Results 457 in Figure 5(a) show that except for the extractability under 0 data sample (the base model knows 458 that it needs to generate Python code, but do not know what codes are needed), the extractability and 459 executability generally improves with the data scale. Results in Figure 5(b) show (1) the base model is unable to generate an effective MAS in zero-shot setting, indicating the necessity for training 460 MAS-GPT. Overall, we observe a promising scaling trend of training MAS-GPT: more data leads to 461 better performance. 462

463 Scaling effects of model size. Here, we compare the performance of MAS-GPT trained based on 464 7B, 14B, and 32B models. Results in Figure 5(c) show that the performance of MAS-GPT improves 465 steadily with the growing model size. Overall, these findings demonstrate the promising potential 466 of MAS-GPT, suggesting that it can be further improved with more diverse, high-quality data and 467 stronger models as the community continues to advance.

468 Case study. To offer an intuitive understanding, we present several examples in Appendix show-casing the query, the MAS-GPT-generated reasoning process, and the MAS-GPT-generated MAS.
470 These show that MAS-GPT can generate query-specific MAS (Section A.1), generalize to unseen queries (Section A.2), generate novel MAS (Section A.3).

472 473

5 CONCLUSION

474

475 Building MAS was time-consuming and resource-intensive. This paper aims to streamline this pro-476 cess into a single LLM inference, making MAS creation as effortless as querying ChatGPT. To this 477 end, we introduce MAS-GPT, an LLM specifically trained to generate executable MAS from arbi-478 trary user queries. Our approach follows a data-driven spirit, leveraging a consistency-oriented data 479 construction pipeline to enhance the coherence and consistency of data pairs. We conduct extensive 480 experiments, comparing MAS-GPT against 10+ baseline methods across 9 benchmarks, using 4 different LLMs as MAS drivers. The results consistently demonstrate that MAS-GPT outperforms all 481 baselines, strongly validating its effectiveness and generalizability. Additionally, we observe MAS-482 GPT's potential to further enhance state-of-the-art reasoning capabilities, as well as its scalability 483 for continued improvements. We believe MAS-GPT can accelerate the adoption of MAS, inspiring 484 future research and real-world applications. 485

486	REFERENCES
487	

507

508

509

513

523

524

525

526

531

532

533

488 Anthropic. Claude 3.5 sonnet. https://www.anthropic.com/news/ claude-3-5-sonnet, 2024. Accessed: 2025-01-22.

- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,
 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared
 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu,
 Yaxi Lu, Yi-Hsin Hung, Chen Qian, et al. Agentverse: Facilitating multi-agent collaboration and
 exploring emergent behaviors. In *The Twelfth International Conference on Learning Representa-*tions, 2024.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
 - Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. In *Forty-first International Conference on Machine Learning*, 2024.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and
 Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine Learning*, pp. 10764–10799. PMLR, 2023.
- Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. Scaling synthetic data creation with 1,000,000 personas. *arXiv preprint arXiv:2406.20094*, 2024.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
 Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021a.
 - Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021b.
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. Metagpt: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2024.
 - Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems. In *NeurIPS 2024 Workshop on Open-World Agents*, 2024.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. T\" ulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for" mind" exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008, 2023.

- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. Encouraging divergent thinking in large language models through multi-agent debate. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pp. 17889–17904, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.992. URL https://aclanthology.org/2024. emnlp-main.992/.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 158–167, 2017.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024a.
- Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=1qvx610Cu7.
- Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. A dynamic llm-powered agent network
 for task-oriented agent collaboration. In *First Conference on Language Modeling*, 2024b.
- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V
 Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective instruction tuning. In *International Conference on Machine Learning*, pp. 22631–22648. PMLR, 2023.
- 567 Cong Lu, Shengran Hu, and Jeff Clune. Intelligent go-explore: Standing on the shoulders of giant
 568 foundation models. In *Automated Reinforcement Learning: Exploring Meta-Learning, AutoML,* 569 and LLMs, 2024.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri
 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
 with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- 574 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

570

576

577 578

579

580

584

- OpenAI. Gpt-40 mini: advancing cost-efficient intelligence. https://openai.com/index/gpt-40-miniadvancing-cost-efficient-intelligence/, 2024a. Accessed: 2025-01-23.
- OpenAI. Introducing openai o1-preview. https://openai.com/index/ introducing-openai-o1-preview/, 2024b. Accessed: 2025-01-22.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *NIPS*, 35:27730–27744, 2022.
- Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen,
 Yusheng Su, Xin Cong, et al. Chatdev: Communicative agents for software development. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15174–15186, 2024a.
- ⁵⁸⁹ Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang,
 ⁵⁹⁰ Zhiyuan Liu, and Maosong Sun. Scaling large-language-model-based multi-agent collaboration.
 ⁵⁹¹ arXiv preprint arXiv:2406.07155, 2024b.
- 593 Qwen. Qwq: Reflect deeply on the boundaries of the unknown. https://qwenlm.github. io/blog/qwq-32b-preview/, 2024. Accessed: 2025-01-27.

- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. arXiv preprint arXiv:2311.12022, 2023.
- 598 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally 599 can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- Kiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark Gerstein. MedAgents: Large language models as collaborators for zero-shot medical reasoning. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 599–621, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.33. URL https://aclanthology.org/2024.findings-acl.33/.
- Kiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R
 Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. In *Forty-first International Conference on Machine Learning*, 2024a.
- Kuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha
 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
 models. In *The Eleventh International Conference on Learning Representations*, 2024b.
- ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 <
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

- Johannes Welbl, Nelson F Liu, and Matt Gardner. Crowdsourcing multiple choice science questions.
 In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pp. 94–106, 2017.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei
 Lin, and Daxin Jiang. Wizardlm: Empowering large pre-trained language models to follow complex instructions. In *The Twelfth International Conference on Learning Representations*, 2024.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xionghui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, et al. Aflow: Automating agentic workflow generation. *arXiv preprint arXiv:2410.10762*, 2024.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyan Luo. LlamaFactory: Unified efficient fine-tuning of 100+ language models. In Yixin Cao, Yang Feng, and Deyi Xiong (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pp. 400–410, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-demos.38. URL https://aclanthology.org/2024.acl-demos.38/.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia
 Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. Advances in Neural Information
 Processing Systems, 36, 2023.

648	Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen
649	Schmidhuber. Gptswarm: Language agents as optimizable graphs. In Forty-first International
650	Conference on Machine Learning, 2024.
651	
652	
653	
654	
655	
656	
657	
658	
659	
660	
661	
662	
663	
664	
665	
666	
667	
668	
670	
671	
672	
672	
674	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	
701	

Table 4: Statistics of MAS-GPT's training dataset. We show the number of samples N_{data} ; the averaged instruction L_{ins} ; the averaged response L_{res} , reasoning L_{rsn} , and MAS length L_{MAS} ; and the number of unique MAS N_{MAS} .

N_{data}	$ L_{ins} $	L_{res}	L_{rsn}	L_{MAS}	N_{MAS}
11442	\mid \sim 75.0 \mid	~ 1062.3	~ 262.5	~ 784.8	7580

Table 5: Case 1: user query and the reasoning process generated by our MAS-GPT. MAS-GPT firstly analyzes the domain, difficulty, and required knowledge, then propose the idea of what multi-agent system is appropriate.

713	Query	The following is a multiple-choice question: Which of the following state-						
714	- •	ments about any two events A and B is true?						
715								
716		The following choices are provided:						
717		(A) $P(A \cup B)$ implies events A and B are independent.						
718		(B) $P(A \cup B) = 0$ implies events A and B are mutually exclusive.						
719		(C) $P(A \cap B)$ implies events A and B are independent.						
720		(D) $P(A \cup B)$ implies events A and B are mutually exclusive. (E) $P(A \cup B)$ implies events A and B are mutually exclusive.						
721		(E) P(A B) = 0 implies events A and B are independent. (E) $P(A \cap B) = 1$ implies events A and B are independent.						
722		(G) P(A + B) = 1 implies events A and B are independent. (G) $P(A + B) = 1$ implies events A and B are independent.						
723		$(H) P(A \cap B) = 1$ implies events A and B are mutually exclusive						
724		(I) $P(A \cup B) = 1$ implies events A and B are mutually exclusive.						
725		(J) $P(A \cap B) = 0$ implies events A and B are independent.						
726								
727		Please finish your answer with "the answer is (X)" where X is the cor-						
728		rect letter choice.						
729	Reasoning	The question is a multiple-choice problem in the domain of probability						
730	generated by	theory, specifically focusing on the relationships between events, such as						
731	MAS-GPT	independence and mutual exclusivity. The difficulty lies in the need for precise						
732		understanding of probability concepts and their implications, as well as the ability to analyze and evaluate multiple statements to identify the correct one						
734		This requires both theoretical knowledge and logical reasoning.						
735								
736		To address this question effectively, a multi-agent system should include						
737		agents that specialize in analyzing the relationships between events in proba-						
738		bility theory. Each agent should focus on a specific aspect of the problem, such						
739		as independence, inducat exclusivity, of the implications of specific probability values (e.g., $P(A \cap B) = 0$ or $P(A \cup B) = 1$). By dividing the task into these						
740		focused analyses, the system ensures comprehensive coverage of the problem						
741		space. Additionally, a final decision-making agent is necessary to synthesize						
742		the insights from all agents and provide a well-reasoned final answer. This						
743		structure ensures that the system can handle the complexity of the question						
744		while leveraging the strengths of specialized reasoning agents. The system is						
745		designed to align with the specific requirements of the question, ensuring that						
746		each agent contributes relevant and precise insights to the solution process.						
747								

A CASE STUDY

748

749

750

702

A.1 MAS-GPT CAN GENERATE QUERY-SPECIFIC MAS 752

From Table 5, we see that given a query, MAS-GPT first provides a reasoning process in analyzing
the domain, difficulty, and required knowledge of the query. Secondly, MAS-GPT analyzes what
multi-agent system is appropriate for addressing the question effectively and decides the required
agents and structure.



```
810
                 Steps:
      11
811
                      1. 5 agents analyze the task independently, each focusing
812
                         on a specific aspect of probability theory.
813
                      2. A final decision-making agent reasons over the
      13
                         solutions and provides the final solution.
814
                  .....
      14
815
      15
816
                  # Step-by-step instruction for each agent to reason and
      16
817
                     generate answer
818
      17
                 instructions = [
819
      18
                      f"Task: {taskInfo}\n\nFocus on analyzing the relationship
                         between events A and B in terms of independence. Solve
820
                          the task.",
821
                      f"Task: {taskInfo}\n\nFocus on analyzing the relationship
      19
822
                         between events A and B in terms of mutual exclusivity.
                          Solve the task.",
823
                      20
824
                          cup B) and P(A \cap B) on the relationship between
825
                          events A and B. Solve the task.",
826
                      f"Task: {taskInfo}\n\nFocus on the definitions and
      21
827
                         properties of probability union and intersection.
828
                         Solve the task.",
                      f"Task: {taskInfo}\n\nFocus on identifying any logical
829
                          inconsistencies or impossibilities in the given
830
                         choices. Solve the task."
831
                 1
      23
832
      24
833
      25
                  # Call the llm to generate each solution
                 cot_results = [self.llm.call_llm(instruction) for instruction
834
      26
                     in instructions]
835
      27
836
                  # Get the instruction for the final decision-making agent
      28
837
                     based on all generated solutions
838
      29
                 final_decision_instruction = self.
                     get_final_decision_instruction(taskInfo, cot_results)
839
      30
840
      31
                  # Call the llm to process the final decision-making
841
                     instruction and generate the final answer
842
                  final_decision_result = self.llm.call_llm(
      32
843
                     final_decision_instruction)
      33
844
                  # Return the final solution
      34
845
                 return final_decision_result
      35
846
      36
847
      37
             def get_final_decision_instruction(self, taskInfo, cot_results):
848
      38
      39
                 Format an instruction for final decision-making based on a
849
                     given task description and a list of solutions.
850
      40
851
      41
                 Args:
852
      42
                      taskInfo (str): A description of the task that needs to be
853
                          completed.
                      cot_results (list): A list containing solutions or
      43
854
                          reasoning steps for the task.
855
      44
856
                 Returns:
      45
857
                      str: A formatted instruction that includes the task
      46
                         description, each solution, and a prompt for final
858
                         decision-making.
859
                  ....
      47
860
      48
861
                  # Initialize the instruction text with a general guideline
      49
862
      50
                 instruction = f"Task:\n{taskInfo}\n\n"
863
      51
                  # Append each solution from cot_results to the instruction
      52
```

for i, result in enumerate(cot_results): instruction += f"Solution {i+1}:\n{result}\n\n" # Number each solution for clarity # Add the final prompt to encourage reasoning over the solutions and provide a final answer instruction += "Given all the above solutions, reason over them carefully and provide a final answer to the task." # Return the complete instruction text return instruction

Table 6: Case 2: user query and the reasoning process generated by our MAS-GPT. MAS-GPT firstly analyzes the domain, required elements, and the difficulty, then propose the idea of what multi-agent system is appropriate.

879 880	Query	An electric dipole consisting of charges of magnitude $1.50nC$ separated by $6.20\mu m$ is in an electric field of strength 1100 N/C . What is the magnitude of
881		the electric dipole moment?
882 883	Reasoning	The question involves calculating the magnitude of the electric dipole moment,
884	generated by	a concept rooted in physics, specifically in the domain of electromagnetism.
885	MAS-01 1	their separation and the resulting dipole moment. While the mathematical
886		computation is straightforward, the problem demands precise application of
887		the formula for the dipole moment and careful handling of units, which adds a
888		layer of complexity.
889		
890		To address this, one would analyze the question by identifying the key
891		variables (charge magnitude, separation distance) and the formula (dipole
892		moment = charge × separation).
893		A multi-agent system is well-suited for this task because it can leverage
894		multiple approaches to ensure accuracy and robustness. First, agents can
895		independently break down the problem into steps, generate Python code to
896		compute the solution, and execute the code to verify the calculations. This en-
897		sures that the solution is not only theoretically correct but also computationally
898		validated. Additionally, other agents can directly generate solutions without
899		code, providing a baseline for comparison. Finally, an evaluation agent can
900		synthesize the outputs from all agents, ensuring consistency and selecting the
901		for scientific problems where both analytical and computational accuracy are
902		critical
903		erneu.

Listing 2: Case 2: Multi-agent system generated by MAS-GPT. MAS-GPT can generalize to unseen queries from SciBench (Wang et al., 2024a), generating an appropriate multi-agent system to handle the query.

```
from utils import *
1
2
   class MAS():
3
       def __init__(self, model_list):
4
           self.llm = LLM(model_list)
5
6
       def forward(self, taskInfo):
7
8
       A multi-agent system for solving math problems by executing code
9
          and directly answering.
10
       Steps:
           1. 2 agents independently solves the problem by breaking it
11
               down into steps and generating code, where each agent
```

organizes the solution based on the code execution results , ensuring clarity and correctness. 2. 2 agents generate a solution directly, which provides baseline solutions especially when code generation is challenging 3. A final agent evaluates all the solutions and determines the final solution. # 4 parallel pipelines to solve the problem independently solutions = [] # The first two pipelines generate code to solve the problem for _ in range(2): answer, output = self.generate_code_get_output(taskInfo) solution = self.organize(taskInfo, answer, output) solutions.append(solution) # The third pipeline generates a solution directly for _ in range(2): solution = self.llm.call_llm(taskInfo) solutions.append(solution) # Determine the final solution based on the generated solutions final_solution = self.get_final_solution(taskInfo, solutions) return final_solution def generate_code_get_output(self, taskInfo): Generate Python code to solve the mathematical problem and execute the code to get the output. Args: taskInfo (str): The mathematical problem to be solved. Returns: a tuple containing: - str: The answer generated by the LLM model. - str: The output of the code execution. code_generation_instruction = f"""You are an expert in solving mathematical problems. **Problem:** {taskInfo} **Instructions:** 1. Analyze the problem and list the steps required to solve it. 2. Generate Python code that can help solve the problem. The code should: - Print important intermediate results in the calculation process, along with clear explanations. - Store the final calculation result in a variable named 'output'. This variable should contain the final result of the computation and be defined at the global scope. - Be directly executable. The code should run and produce a result when executed. Wrap your final code solution in <Code Solution> and </Code Solution>. For example: <Code Solution> Your function code here </Code Solution> # Call 'generate_and_extract_code' to generate answer and extract the code answer, code = generate_and_extract_code(llm=self.llm, prompt= code_generation_instruction)

```
972
                  # Call 'execute_code' to execute the generated code and get
      60
973
                      output
974
                  output = execute_code(code)
      61
                  return answer, output
975
      62
      63
976
              def organize(self, taskInfo, answer, result):
      64
977
      65
978
              Organize the solution based on the code execution results.
      66
979
              Args:
      67
980
      68
                  taskInfo (str): The mathematical problem to be solved.
                  answer (str): The initial solution generated by the LLM model.
      69
981
                  result (str): The output of the code execution.
      70
982
              Returns:
      71
983
                  str: The organized solution based on the code execution
      72
984
                      results.
              .....
985
      73
                  organize_instruction = f"""**Problem:**
      74
986
          {taskInfo}
      75
987
      76
          **Initial Solution:**
988
      77
          {answer}
989
          **Code Execution Result:**
      78
990
          {result}
      79
         To solve the **Problem**, the **Initial Solution** provides steps and
      80
991
             python code for calculations. The **Code Execution Result** is the
992
               output of the code.
993
         Based on the **Initial Solution** and **Code Execution Result**,
      81
994
              provide a final solution to the problem. Include the results of
995
              the code calculation in your response. Your final response should
              be complete as if you are directly answering the problem."""
996
                  solution = self.llm.call_llm(organize_instruction)
      82
997
      83
                  return solution
998
      84
999
              def get_final_solution(self, taskInfo, solutions):
      85
              ....
1000
      86
              Evaluate the solutions provided by the agents and determine the
      87
1001
                  final solution.
1002
      88
              Args:
1003
                  taskInfo (str): The mathematical problem to be solved.
      89
1004
                  solutions (list): A list of solutions provided by the agents.
      90
              Returns:
1005
      91
                  str: The final solution to the mathematical problem.
      92
1006
              .....
      93
1007
                  # Prepare the solutions in a numbered format for evaluation
      94
1008
                  solutions_str = ""
      95
1009
                  for i, solution in enumerate(solutions):
      96
                      solutions_str += f"Solution {i+1}\n{solution}\n\n"
1010
      97
      98
1011
                  final_solution_instruction = f"""**Problem:**
      99
1012
          {taskInfo}
      100
1013
          **Solutions:**
      101
1014
      102
          {solutions_str}
          Several solutions are given to solve the problem. Carefully evaluate
1015
      103
              these solutions. This consistency in answers is crucial for
1016
              determining the most reliable solution.
1017
      104
          You should provide a final solution to the problem. Your final
1018
              response should be complete as if you are directly answering the
             problem."""
1019
                  final_solution = self.llm.call_llm(final_solution_instruction)
      105
1020
                  return final_solution
      106
1021
1022
1023
1024
```

simultaneously as in the Bertrand model. Compute the prices in the na equilibrium. The following choices are provided: (A) 30 (B) 50 (C) 0 (D) 25 (E) 150 (F) 75 (G) 90 (H) 120 (I) 100 (I) 60 Please finish your answer with "the answer is (X)" where X is the correct let choice. Reasoning generated by MAS-GPT MAS-GPT The question pertains to a classic economic problem in the domain of gan theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interacti between two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive t equilibrium price. To address such a question, one must first break it down into its cc components; the demand curve, the firms' identical products, and the simul neous price-setting behavior. A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By emplé ing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a rob exploration of the problem space. This approach is particularly valuable f a question like this, where multiple solution paths (e.g., algebraic derivati economic intuition) can lead to the correct answer. The iterative refinem process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solutio ensuring that the final answer is both logically sound and consistent with t principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome.	Query	The following is a multiple-choice question: Suppose that there are two firm in the market facing no costs of production and a demand curve given by $C = 150 - P$ for their identical products. Suppose the two firms choose price
The following choices are provided:(A) 30(B) 50(C) 0(D) 25(E) 150(F) 75(G) 90(H) 120(I) 100(J) 60MAS-GPTThe question pertains to a classic economic problem in the domain of gat theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interacti between two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive t equilibrium price.To address such a question, one must first break it down into its cc components: the demand curve, the firms' identical products, and the simul neous price-setting behavior.A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By emplei ing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robi exploration of the problem space. This approach is particularly valuable f a question like this, where multiple solution paths (e.g., algebraic derivati economic intuition) can lead to the correct answer. The iterative refinem process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solutio ensuring that the final answer is both logically sound and consistent with the principles of the question and ensures a reliable outcome.		simultaneously as in the Bertrand model. Compute the prices in the nas equilibrium.
 (A) 30 (B) 50 (C) 0 (D) 25 (E) 150 (F) 75 (G) 90 (H) 120 (I) 100 (J) 60 Please finish your answer with "the answer is (X)" where X is the correct let choice. Reasoning generated by MAS-GPT The question pertains to a classic economic problem in the domain of gan generated by mAS-GPT The question pertains to a classic economic problem in the domain of gan generated by mAS-GPT The question pertains to a classic economic problem in the domain of gan generated by microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interacti between two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive the equilibrium price. To address such a question, one must first break it down into its concomponents: the demand curve, the firms' identical products, and the simulaneous price-setting behavior. A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By empleting agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robic exploration of the problem space. This approach is particularly valuable ta question like this, where multiple solution paths (e.g., algebraic derivatic economic intuition) can lead to the correct answer. The iterative refineme process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a r		The following choices are provided:
 (B) 50 (C) 0 (D) 25 (E) 150 (F) 75 (G) 90 (H) 120 (I) 100 (J) 60 Please finish your answer with "the answer is (X)" where X is the correct let choice. Reasoning generated by MAS-GPT The question pertains to a classic economic problem in the domain of gar theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interacti between two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive the equilibrium price. To address such a question, one must first break it down into its concomponents: the demand curve, the firms' identical products, and the simulaneous price-setting behavior. A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By empting agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robit exploration of the problem space. This approach is particularly valuable to a question like this, where multiple solution paths (e.g., algebraic derivatic economic intuition) can lead to the correct answer. The iterative refinemend process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome. 		(A) 30
 (C) 0 (D) 25 (E) 150 (F) 75 (G) 90 (H) 120 (I) 100 (J) 60 Please finish your answer with "the answer is (X)" where X is the correct let choice. Reasoning generated by MAS-GPT The question pertains to a classic economic problem in the domain of gat theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interacti between two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive t equilibrium price. To address such a question, one must first break it down into its cc components: the demand curve, the firms' identical products, and the simul neous price-setting behavior. A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By emploing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robe exploration of the problem space. This approach is particularly valuable ta question like this, where multiple solution paths (e.g., algebraic derivatid economic intuition) can lead to the correct answer. The iterative refineme process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome. 		(B) 50
 (E) 150 (E) 150 (F) 75 (G) 90 (H) 120 (I) 100 (J) 60 Please finish your answer with "the answer is (X)" where X is the correct let choice. Reasoning generated by MAS-GPT The question pertains to a classic economic problem in the domain of gat theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interactive between two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive the equilibrium price. To address such a question, one must first break it down into its cc components: the demand curve, the firms' identical products, and the simul neous price-setting behavior. A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By emploing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robe exploration of the problem space. This approach is particularly valuable if a question like this, where multiple solution paths (e.g., algebraic derivatid economic intuition) can lead to the correct answer. The iterative refineme process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome. 		(C) 0 (D) 25
 (F) 75 (G) 90 (H) 120 (J) 100 (J) 60 Please finish your answer with "the answer is (X)" where X is the correct let choice. Reasoning generated by MAS-GPT The question pertains to a classic economic problem in the domain of gat theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interactibetween two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive the equilibrium price. To address such a question, one must first break it down into its concomponents: the demand curve, the firms' identical products, and the simul neous price-setting behavior. A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By emplain agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robe exploration of the problem space. This approach is particularly valuable a question like this, where multiple solution paths (e.g., algebraic derivative conomic intuition) can lead to the correct answer. The iterative efficient of process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome. 		(E) 25 (E) 150
(G) 90 (H) 120 (I) 100 (J) 60Please finish your answer with "the answer is (X)" where X is the correct let choice.Reasoning generated by MAS-GPTThe question pertains to a classic economic problem in the domain of gat theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interacti between two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive t equilibrium price.To address such a question, one must first break it down into its co components: the demand curve, the firms' identical products, and the simul neous price-setting behavior.A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By empli ing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a rob exploration of the problem space. This approach is particularly valuable a question like this, where multiple solution paths (e.g., algebraic derivati economic intuition) can lead to the correct answer. The iterative efficient process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with t principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome.		(F) 75
 (H) 120 (I) 100 (J) 60 Please finish your answer with "the answer is (X)" where X is the correct let choice. Reasoning generated by MAS-GPT The question pertains to a classic economic problem in the domain of gat theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interactive tequilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive tequilibrium price. To address such a question, one must first break it down into its concomponents: the demand curve, the firms' identical products, and the simul neous price-setting behavior. A multi-agent system is well-suited for this task because it can leveral multiple perspectives to ensure accuracy and depth in reasoning. By emploing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robic exploration of the problem space. This approach is particularly valuable for a question like this, where multiple solution paths (e.g., algebraic derivatio economic intuition) can lead to the correct answer. The iterative refinemer process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome. 		(G) 90
 (1) 100 (1) 60 Please finish your answer with "the answer is (X)" where X is the correct let choice. Reasoning generated by MAS-GPT The question pertains to a classic economic problem in the domain of gat theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interactive tequilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive tequilibrium price. To address such a question, one must first break it down into its concomponents: the demand curve, the firms' identical products, and the simul neous price-setting behavior. A multi-agent system is well-suited for this task because it can leveral multiple perspectives to ensure accuracy and depth in reasoning. By emploing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robic exploration of the problem space. This approach is particularly valuable of a question like this, where multiple solution paths (e.g., algebraic derivatic economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome. 		(H) 120 (I) 100
Please finish your answer with "the answer is (X)" where X is the correct let choice.Reasoning generated by MAS-GPTThe question pertains to a classic economic problem in the domain of gar 		(J) 60
Please finish your answer with "the answer is (X)" where X is the correct let choice.Reasoning generated by MAS-GPTThe question pertains to a classic economic problem in the domain of gar theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interacti between two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive t equilibrium price.To address such a question, one must first break it down into its co components: the demand curve, the firms' identical products, and the simul neous price-setting behavior.A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By emplo ing agents that independently analyze the problem, generate solutions, a a refine their responses through iterative feedback, the system ensures a rob exploration of the problem space. This approach is particularly valuable t a question like this, where multiple solution paths (e.g., algebraic derivatio economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with t principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome.		
Please finish your answer with "the answer is (X)" where X is the correct let choice.Reasoning generated by MAS-GPTThe question pertains to a classic economic problem in the domain of gar theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interacti between two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive t equilibrium price.To address such a question, one must first break it down into its cc components: the demand curve, the firms' identical products, and the simul neous price-setting behavior.A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By emploi ing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robi exploration of the problem space. This approach is particularly valuable ta a question like this, where multiple solution paths (e.g., algebraic derivati conomic intuition) can lead to the correct answer. The iterative refineme process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with t principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome.		
Reasoning generated by MAS-GPTThe question pertains to a classic economic problem in the domain of gar theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interact ibetween two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive t equilibrium price.To address such a question, one must first break it down into its cc components: the demand curve, the firms' identical products, and the simul neous price-setting behavior.A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By emploi ing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robi exploration of the problem space. This approach is particularly valuable d a question like this, where multiple solution paths (e.g., algebraic derivati economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with t principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome.		Please finish your answer with "the answer is (X) " where X is the correct letter shoice
 Reasoning generated by MAS-GPT MAS-GPT The question pertains to a classic economic problem in the domain of gat theory and microeconomics, specifically focusing on the Bertrand model price competition. The difficulty lies in understanding the strategic interactive between two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive tequilibrium price. To address such a question, one must first break it down into its concomponents: the demand curve, the firms' identical products, and the simul neous price-setting behavior. A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By emploing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robit exploration of the problem space. This approach is particularly valuable for a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome. 	Dessering	The spectrum partning to a classic economic method in the density of econ
 MAS-GPT massing the structure of the st	generated by	theory and microeconomics, specifically focusing on the Bertrand model of
 between two firms, the implications of no production costs, and the Na equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive t equilibrium price. To address such a question, one must first break it down into its co components: the demand curve, the firms' identical products, and the simul neous price-setting behavior. A multi-agent system is well-suited for this task because it can levera multiple perspectives to ensure accuracy and depth in reasoning. By emploing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robe exploration of the problem space. This approach is particularly valuable to a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome. 	MAS-GPT	price competition. The difficulty lies in understanding the strategic interactio
 equilibrium concept. Solving this requires both theoretical knowledge economic models and the ability to apply mathematical reasoning to derive t equilibrium price. To address such a question, one must first break it down into its cc components: the demand curve, the firms' identical products, and the simul neous price-setting behavior. A multi-agent system is well-suited for this task because it can leveral multiple perspectives to ensure accuracy and depth in reasoning. By emploing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robut exploration of the problem space. This approach is particularly valuable a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome. 		between two firms, the implications of no production costs, and the Nas
 To address such a question, one must first break it down into its co components: the demand curve, the firms' identical products, and the simul neous price-setting behavior. A multi-agent system is well-suited for this task because it can leveral multiple perspectives to ensure accuracy and depth in reasoning. By emploing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robe exploration of the problem space. This approach is particularly valuable for a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome. 		equilibrium concept. Solving this requires both theoretical knowledge of economic models and the ability to apply mathematical reasoning to derive the
To address such a question, one must first break it down into its concomponents: the demand curve, the firms' identical products, and the simulation neous price-setting behavior. A multi-agent system is well-suited for this task because it can leveral multiple perspectives to ensure accuracy and depth in reasoning. By emploing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robe exploration of the problem space. This approach is particularly valuable for a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome.		equilibrium price.
 components: the demand curve, the firms' identical products, and the simul neous price-setting behavior. A multi-agent system is well-suited for this task because it can leveral multiple perspectives to ensure accuracy and depth in reasoning. By emploing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robe exploration of the problem space. This approach is particularly valuable for a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome. 		To address such a question, one must first break it down into its cor
A multi-agent system is well-suited for this task because it can leveral multiple perspectives to ensure accuracy and depth in reasoning. By emplo- ing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robe exploration of the problem space. This approach is particularly valuable f a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with t principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome.		components: the demand curve, the firms' identical products, and the simulta
A multi-agent system is well-suited for this task because it can leveral multiple perspectives to ensure accuracy and depth in reasoning. By emplo- ing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robe exploration of the problem space. This approach is particularly valuable to a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome.		neous price-setting behavior.
multiple perspectives to ensure accuracy and depth in reasoning. By emploing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robut exploration of the problem space. This approach is particularly valuable for a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome.		A multi-agent system is well-suited for this task because it can leverage
ing agents that independently analyze the problem, generate solutions, a refine their responses through iterative feedback, the system ensures a robu exploration of the problem space. This approach is particularly valuable f a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome.		multiple perspectives to ensure accuracy and depth in reasoning. By employ
exploration of the problem space. This approach is particularly valuable is a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome.		ing agents that independently analyze the problem, generate solutions, an
a question like this, where multiple solution paths (e.g., algebraic derivation economic intuition) can lead to the correct answer. The iterative refinement process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with the principles of the Bertrand model. This structured reasoning process align with the complexity of the question and ensures a reliable outcome.		refine their responses through iterative feedback, the system ensures a robut exploration of the problem space. This approach is particularly valuable for
economic intuition) can lead to the correct answer. The iterative refinemed process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with t principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome.		a question like this, where multiple solution paths (e.g., algebraic derivation
process allows for cross-verification of solutions, reducing the likelihood errors. Finally, a decision-making agent synthesizes the diverse solution ensuring that the final answer is both logically sound and consistent with t principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome.		economic intuition) can lead to the correct answer. The iterative refinement
ensuring that the final answer is both logically sound and consistent with t principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome.		process allows for cross-verification of solutions, reducing the likelihood of armore Einstein a desirier making egent surtherings the diverse solution
principles of the Bertrand model. This structured reasoning process alig with the complexity of the question and ensures a reliable outcome.		ensuring that the final answer is both logically sound and consistent with the
with the complexity of the question and ensures a reliable outcome.		principles of the Bertrand model. This structured reasoning process align
		with the complexity of the question and ensures a reliable outcome.

1080	6	
1081	0 7	<pre>def forward(self, taskInfo):</pre>
1082	8	
1083	9	A multi-agent system for solving general tasks.
1084	10	
1085	11	Steps:
1086	12	2 Each agent reflects on the solutions and provides an
1087	15	improved solution.
1088	14	3. A final decision-making agent reasons over the improved
1089		solutions and provides the final solution.
1090	15	
1091	16	# Step-by-step instruction for each agent to reason and
1092	17	generate answer instruction = f"Task· {taskInfo}\n\nPlease solve the task "
1093	18	
1094	19	<pre># Set the number of solutions to generate; using 5 for variety</pre>
1095		and diversity
1096	20	N = 5
1097	21	# Call the lim to generate each solution
1098	22)]
1099	23	
1100	24	# Get the instruction for the self-refine process based on all
1101		generated solutions
1102	25	<pre>self_refine_instruction = self.get_self_refine_instruction(</pre>
1102	26	taskinio, cot_results)
1104	20	# Call the llm to refine each solution
1104	28	refined_results = [self.llm.call_llm(self_refine_instruction)
1105		<pre>for _ in range(N)]</pre>
1100	29	
1107	30	# Get the final decision-making instruction based on all
1100	31	final decision instruction = self.
11109		<pre>get_final_decision_instruction(taskInfo, refined_results)</pre>
1110	32	
1110	33	# Call the llm to process the final decision-making
1112	24	instruction and generate the final answer
1113	34	final decision instruction)
1114	35	
1115	36	# Return the final solution
1110	37	return final_decision_result
1117	38	
1118	39 40	<pre>def get_self_refine_instruction(self, taskinfo, cot_results): """</pre>
1119	40	Format an instruction for self-refinement based on a given
1120		task description and a list of solutions.
1121	42	
1122	43	Args:
1123	44	taskinio (str): A description of the task that needs to be
1124	45	cot results (list): A list containing solutions or
1125	15	reasoning steps for the task.
1126	46	
1127	47	Returns:
1128	48	str: A formatted instruction that includes the task
1129		description, each solution, and a prompt for self-
1130	49	IIII
1131	50	
1132	51	# Initialize the instruction text with a general guideline
1133	52	<pre>instruction = f"Task:\n{taskInfo}\n\n"</pre>
	53	

Append each solution from cot_results to the instruction for i, result in enumerate(cot_results): instruction += f"Solution {i+1}:\n{result}\n\n" # Number each solution for clarity # Add the final prompt to encourage self-refinement and improvement instruction += "Given all the above solutions, reason over them carefully and provide an improved solution to the task." # Return the complete instruction text return instruction def get_final_decision_instruction(self, taskInfo, refined_results): Format an instruction for final decision-making based on a given task description and a list of refined solutions. Args: taskInfo (str): A description of the task that needs to be completed. refined_results (list): A list containing refined solutions for the task. Returns: str: A formatted instruction that includes the task description, each refined solution, and a prompt for final decision-making. # Initialize the instruction text with a general guideline instruction = f"Task:\n{taskInfo}\n\n" # Append each refined solution from refined_results to the instruction for i, result in enumerate(refined_results): instruction += f"Solution {i+1}:\n{result}\n\n" # Number each solution for clarity # Add the final prompt to encourage reasoning over the solutions and provide a final answer instruction += "Given all the above solutions, reason over them carefully and provide a final answer to the task." # Return the complete instruction text return instruction VISUALIZATION OF OUR MAS POOL В We visualize several MAS in our MAS pool in Figure 7. С ADDITIONAL EXPERIMENTAL SETUPS. C.1 DESCRIPTIONS OF DATASETS

We provide an overall descriptions of the training and testing datasets in Table 9.



242 243	Method	HUMANEVAL	HUMANEVAL-PLUS	MATH
244	CHATDEV (QIAN ET AL., 2024A)	83.33	84.04	62.07
	MAS-GPT (OURS)	91.18	87.23	77.59

1246 1247 1248

Table 8: Comparisons with task-specific method, ChatDev (Qian et al., 2024a), which is specifically designed for software development.

Purpos	e Dataset Name	Domain	Sub-Domains	Sample Number
Trainin	MATH (Hendrycks et al., 2021b)	Math	Counting & Probabil- ity Geometry Algebra Number Theory Precalculus Prealgebra	6000
	GSM8K (Cobbe et al., 2021)	Math	-	1000
	GSM-Hard (Gao et al., 2023)	Math	-	319
	AQUA-RAT (Ling et al., 2017)	Reasoning	-	1000
	MBPP (Austin et al., 2021)	Code	-	374
	SciQ (Welbl et al., 2017)	General QA	-	2000
			Humanities	
	MMLU (Hendrycks et al., 2021a)	General OA	Social Science	1529
			STEM	
			Others	
	MATH (Hendrycks et al., 2021b)	Math	Counting & Probabil- ity Geometry Algebra Number Theory Precalculus Prealgebra	500
Testing		Ma	Intermediate Algebra	500
	GSM8K (Cobbe et al., 2021)	Math	-	500
	HumanEval (Chan et al. 2021)	Code	-	300 164
	HumanEval Plus (Line et al., 2021)	Code	-	164
	GPOA (Rein et al. 2023)	Science	-	448
	SciBench (Wang et al. 2023)	Science	-	500
	Serbenen (Trang et al., 2024a)	berenee	Humanities	200
			Social Science	-
	MMLU (Hendrycks et al., 2021a)	General QA	STEM	500
			Others	
	AIME-2024	Math	-	30
	1	1		

1282 1283 1284

Table 9: Descriptions of benchmarks

1285 C.2 EVALUATIONS 1286

In this section, we detail our evaluation approach. For queries with ground truth answers, we employ 1287 LLMs to extract the MAS output and compare it with the ground truth. For code benchmarks like 1288 HumanEval and MBPP, we assess correctness using test cases. 1289

1290 LLM-based Evaluation with Ground-Truth Answer We utilize LLMs to perform evaluation with 1291 ground-truth answer. However, direct evaluation against the ground truth is incompatible as the LLM 1292 annotates the response itself. To address this issue, we adopt a two-step evaluation process based 1293 on the prompts used in AutoGen (Wu et al., 2023), first extract the answer, then evaluation. Here the responses generated by multi-agent systems (MAS) are often unstructured and irregular, making 1294 it difficult to extract the final answer to a query using rule-based methods. To avoid extraction 1295 errors that could impact the evaluation of MAS performance, we use LLMs to automate the answer

1296 Table 10: The prompt for intra-consistency-oriented pair refinement. This prompt is fed to GPT-4o-1297 2024-11-20 to adjust the MAS and generate a reasoning statement. The prompt is integrated with a 1298 user query and the selected MAS represented by Python code. 1299 1300 I will give you a question and a multi-agent system. The multi-agent system is described in 1301 the format of Python code, where each agent is represented by an agent-specific instruction 1302 and one call_llm. Though the multi-agent system can answer the question, it may not be the 1303 best one. You task is return me two things: an improved multi-agent system and a paragraph. 1304 1305 The improved multi-agent system should be more related to the question, while basically, 1306 try not to change the architecture compared to the original multi-agent system. - For example, if the multi-agent system is in a parallel structure (e.g., 5 parallel agents generate answer and 1 agent select the best answer), you may keep the structure unchanged 1309 while only changing each parallel agent's instruction. 1310 - If the multi-agent system is already suitable, you may only modify the instructions in the multi-agent system more relevant to the question while leaving the structure unchanged. 1311 - If you think additional agents are required (e.g., the question is difficult and complex), you 1312 may add some related expert agents to enhance the multi-agent system. 1313 1314 The paragraph should first analyze the question itself, from the perspectives of domain and 1315 difficulty. Then, you should provide a reasoning process to bridge the question and the 1316 improved multi-agent system. The reasoning process should be in the views of that how 1317 one analyzes the question and objectively thinks about what multi-agent system is needed. 1318 Then, the reasoning process can finally and logically lead to the improved multi-agent 1319 system. Do not mention "this multi-agent system", or "the improved multi-agent system", 1320 rather, say "a multi-agent system" instead. Do not mention the original multi-agent system 1321 or the original structure. 1322 Please follow the following format requirements: 1323 - The improved multi-agent system should be included between <CODE> and </CODE> 1324 - The paragraph should be included between <PARAGRAPH> and </PARAGRAPH> 1326 Please firstly generate the multi-agent system and then generate the paragraph. The 1327 paragraph should analyze about the question and the generated multi-agent sytem, such that 1328 when one sees the (question, paragraph, the improved multi-agent system) triplet (wihtout the original multi-agent system), one can understand the reasoning process behind the 1330 improved multi-agent system. Notice! The paragraph should never mention the original 1331 multi-agent system or the original structure. 1332 The question is: 1333 {query} 1334 1335 The multi-agent system is: 1336 {MAS code} 1337 1338 1339 1340 extraction process. Specifically, we prompt the LLM to extract the answer from the MAS response 1341 based on predefined rules and then ask the LLM to compare it with the ground truth. The prompts 1342 used for this process are detailed in Table 12.

Code Evaluation with Test Cases We evaluate the MAS performance on coding tasks based on pass rate on test cases, with a two step approach: first, prompting the LLM to extract the code from the MAS response, and second, executing it in a coding environment to calculate the pass rate; see the prompts used for extract code and functions in Table 13.

- 1348
- 1349

1364Table 11: The prompt for generating a reasoning process if the refined MAS fails. This prompt is1365fed to GPT-4o-2024-11-20 to generate a reasoning statement. The prompt is integrated with a user1366query and a selected MAS represented by Python code.

I will give you a question and a multi-agent system. The multi-agent system is described in the format of Python code, where each agent is represented by an agent-specific instruction and one call_llm. You task is return me a paragraph.

The paragraph should first analyze the question itself, from the perspectives of domain and difficulty. Then, you should provide a reasoning process to bridge the question and the provided multi-agent system. The reasoning process should be in the angle of views that how one analyzes the question and objectively thinks about what multi-agent system is needed. Then, the reasoning process can finally and logically lead to the provided multi-agent system. Do not mention "this multi-agent system", or "the provided multi-agent system", rather, say "a multi-agent system" instead.

The paragraph should analyze about the question and the provided multi-agent system, such that when one sees the (question, paragraph, the provided multi-agent system) triplet, one can understand the reasoning process behind the provided multi-agent system.

Remember, the paragraph should be included between $;PARAGRAPH_{i}$ and $;/PARA-GRAPH_{i}$.

The question is: {query}

The provided multi-agent system is: {MAS code}

1404 1405 1406 1407 1408 Table 12: Prompts for extract answer and answer evaluation. 1409 1410 You are a helpful AI assistant tasked with extracting the final answer from a provided 1411 solution. 1412 **Input:** 1413 1. A problem statement, prefixed with "===Problem: <problem>". 1414 2. A solution to the problem, prefixed with "===Solution:<solution>". 1415 1416 **Problem and Solution:** 1417 ===Problem: {query} 1418 1419 ===Solution: {response} 1420 1421 **Instructions:** 1422 - Carefully analyze the solution and extract the final answer in reply: "The answer is 1423 <answer extracted> in reply". - If the solution does not contain a final answer (e.g., only reasoning, code without execu-1424 tion, or incomplete information), respond with: "The reply doesn't contain an answer." 1425 - Ensure that the extracted answer is exactly as presented in the solution. Do not infer or use 1426 external knowledge. Do not execute the code yourself. 1427 - Remember, Never execute the code yourself! Never doing any computation yourself! Just 1428 extract and output the existing answer! 1429 1430 1431 1432 You are a helpful AI assistant. You will use your coding and language skills to verify the 1433 answer. 1434 You are given: 1. A problem, which is going to start like "===Problem: <problem>". 1435 2. A ground truth answer, which is going to start like "===Ground truth answer:". 1436 3. A reply with the answer to the problem, which are going to start like "===Reply:". 1437 Please do the following: 1438 1. Extract the answer in reply: "The answer is <answer extracted> in reply". 1439 2. Check whether the answer in reply matches the ground truth answer. When comparison is 1440 not obvious (for example, 3*sqrt(6) and 7.348), you may compare by calculation, allowing 1441 a small margin of error. 1442 3. After everything is done, please give each reply a comment like the following options: 1443 - "The answer is correct." - "The answer is approximated but should be correct. Correct Answer: <ground truth 1444 answer> | Answer extracted: <answer extracted>." 1445 - "The answer is incorrect. Correct Answer: <ground truth answer> | Answer 1446 extracted: <answer extracted>." 1447 - "The reply doesn't contain an answer." 1448 Here are the problem, the ground truth answer and the reply: 1449 ===Problem: {query} 1450 1451 ===Ground truth answer: {ground_truth_answer} 1452 1453 ===Reply: {Reply} 1454 1455 1456

1458	
1459	
1460	
1/61	
1401	
1462	
1463	
1464	
1465	
1466	
1/67	
1407	
1468	
1469	
1470	
1471	
1472	
1/172	
1473	
14/4	
1475	
1476	
1477	Table 13. Prompts for extract code and functions
1478	
1/170	
1475	You are given a **Problem** and a **Solution**. The **Problem** asks for a code
1480	function. Extract the final code function from the **Solution**.
1481	**Problem:**
1482	{query}
1483	
1484	**Solution:**
1/05	{solution}
1400	
1486	Please follow the following rules:
1487	Only output the code function that exists in the **Solution**, without any additional ex
1488	- Only output the code function that exists in the "Solution", without any additional ex-
1489	planation of content.
1490	- Do not modify any part of the code function.
1/01	- Remove parts like 'example use' or 'test cases'.
1491	- If the **Solution** does not contain a code function, respond with: "The reply doesn't
1492	contain a code function."
1493	
1494	
1495	
1496	
1/107	
1/00	
1490	
1499	
1500	
1501	
1502	
1503	
1505	
1004	
1505	
1506	
1507	
1508	
1500	
1509	
1509 1510	
1509 1510 1511	