
BEYOND ZERO-SHOT: A HIERARCHICAL APPROACH TO MULTI-AGENT COLLABORATION FOR IMPROVED LLM ACCURACY

Giovanpaolo Vrenna

Giovanpaolo.vrenna@keio.jp

Shigeki Noguchi

noguchi_shigeki@keio.jp

Naoki Watson

naokiwatson@keio.jp

Keita Fujie

fk77663@keio.jp

ABSTRACT

When Large Language Model agents engage in debate, they collectively enhance their intelligence. This study explores the impact of hierarchical communication structures on multi-agent debate systems, introducing two distinct hierarchical topologies along with a method for their implementation. Our results indicate that these structures improve reasoning performance, achieving approximately 20% higher accuracy compared to zero-shot agents. Additionally, we conducted a comparative analysis of shallow and deep hierarchical topologies. While our findings do not explicitly reveal a performance difference, our analysis of conversation history suggests that deeper hierarchical structures introduce greater redundancy in the reasoning process, which may, in turn, enhance self-correction

1 Introduction

Since the introduction of the transformer architecture in 2017[1] and the release of OpenAI’s GPT models[2][3][4], the machine learning community has primarily focused on transformer-based models. The discovery of scaling laws[5] further shifted research priorities from optimizing model efficiency to increasing model and dataset sizes. As models were trained on vast datasets sourced from the internet, scaling laws initially continued to hold, but their limitations have since been questioned. While it remains uncertain whether these laws have truly reached a ceiling, it is evident that most publicly available raw text data has already been utilized. Consequently, research efforts begun shifting toward improving model inference quality. As of February 2025, DeepSeek R1, one of the leading models known for its strong reasoning capabilities, had been further refined through

reinforcement learning[6]. Post-training techniques for foundational models have become vital for achieving high reasoning performance.

Moreover, recent research demonstrates that LLMs can achieve enhanced performance through techniques like prompt engineering, without updating model weights. For example, Chain of Thought (CoT) reasoning[7] improves an LLM’s ability to solve complex problems by simply instructing it to “think step by step.” Similarly, reflection prompting[8][9] models to think about their own thinking before finalizing an answer—has proven effective in refining their reasoning capabilities. Moreover, employing a multi-agent debate strategy, where several LLM agents engage in structured discussions, has been shown to further boost performance compared to a single-agent approach [10].

In this study, we analyzed the performance of different multi-agent debate system hierarchies. Specifically, we compared the effectiveness of a zero-shot agent against various multi-agent configurations, assessing their impact on accuracy, computational cost, and time efficiency. Our findings provide insights into the trade-offs associated with different debate structures and highlight the advantages of multi-agent collaboration in enhancing reasoning capabilities. Furthermore, we discuss adaptive hierarchical structures that dynamically evolve based on performance metrics for more accurate multi-agent reasoning frameworks.

2 Background

Collective intelligence refers to the phenomenon where the combined capabilities of multiple individuals or agents—whether human or machine—yield problem-solving and decision-making outcomes that surpass those of any single participant. In the context of large language models, Du et al(2023)[10] explored their collective performance within a multi-agent debate system. The debate workflow follows these steps:

1. Given a question, multiple agents generate independent answers.
2. Agents review each other’s responses and reasoning.
3. Each agent revises its response based on feedback from others.
4. This iterative process continues for multiple rounds.

The communication topology in the study was fully connected, meaning each agent had access to all other agents’ responses at every debate round. There was no hierarchical structure, and all agents played an equal role in refining their answers. Instead of fixed pairwise debates, responses were broadcast to all agents, allowing them to critique and update their reasoning iteratively. This setup facilitated iterative consensus formation, where incorrect or uncertain answers were gradually

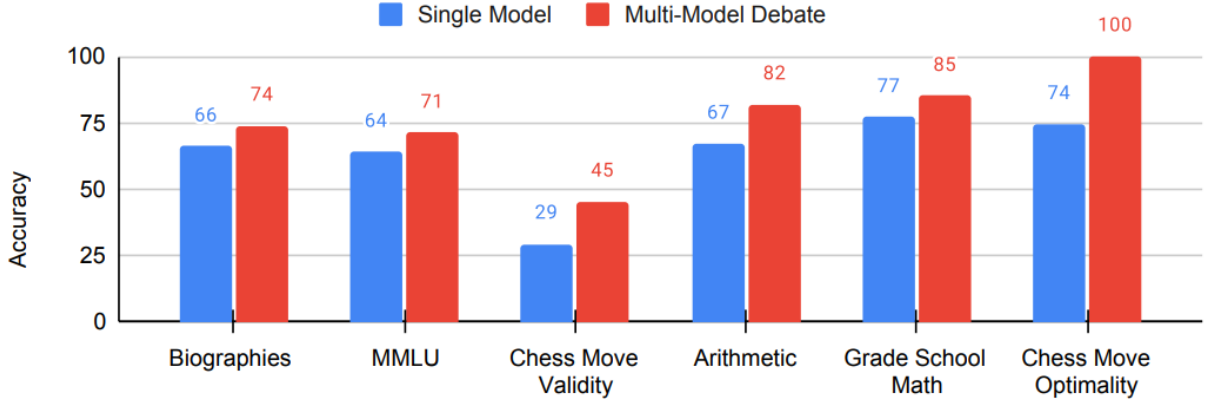


Figure 1: Multiagent Debate Improves Reasoning and Factual Accuracy

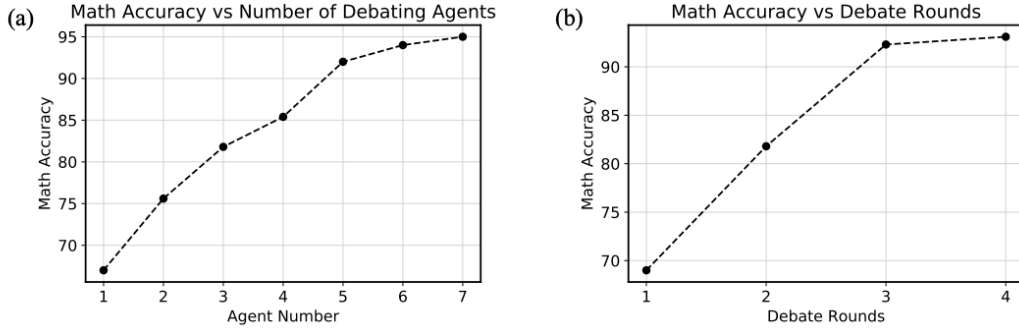


Figure 2: (a) Performance with Increased Agents. (b) Performance with Increased Rounds

refined. The debate process did not rely on a central judge or leader; instead, all agents contributed equally.

The study evaluated the multi-agent debate framework across six benchmarks covering reasoning and factual accuracy tasks. For reasoning, they tested arithmetic expressions, grade school math (GSM8K), and chess move prediction, measuring accuracy and logical consistency. For factual accuracy, they assessed biography generation, MMLU (Massive Multitask Language Understanding), and chess move validity, testing the correctness of generated facts. The multi-agent debate consistently outperformed single-agent approaches, self-reflection methods, and majority voting as summarized in Figure 1.

Although this result is surprising, the performance of the multi-agent debate system plateaus as the number of agents or rounds increases, eventually converging to a certain level (see Figure 2). Based on these observations, we hypothesize that incorporating dynamic, hierarchical communication topologies could further enhance the system. To implement the debate system and its communication topology, we leveraged the Autogen framework[11], which provides the flexibility to define both the topology and the agents’ behavior.

3 Implementation

3.1 Framework: Microsoft AutoGen

AutoGen is a multi-agent framework developed by Microsoft. It offers natively a more flexible agent topology, enhanced LLM inference, and dynamic agent collaboration. These features enable us to build a more robust and flexible multi-agent system [11].

These are the key components of AutoGen that we used in our project:

- **ConversableAgent:** In AutoGen’s framework, the `ConversableAgent` class serves as a customizable foundation for agents capable of engaging in conversations with other agents, humans, and tools to accomplish tasks collaboratively.
- **GroupChat:** A `GroupChat` is a collection of agents that can communicate with each other.
- **GroupChatManager:** In AutoGen’s multi-agent group chat, the `GroupChatManager` plays a major role in facilitating agent communication. When an agent generates a response, the `GroupChatManager` broadcasts it to all participating agents in `GroupChat`. This broadcasting mechanism ensures that all agents remain informed of the ongoing conversation, allowing them to contribute effectively to collaborative tasks. The `GroupChatManager` also manages the flow of the conversation by selecting the next speaker, by orchestrating a coherent and organized dialogue among the agents.
- **SocietyOfMindAgent:** The `Society of Mind Agent` in Microsoft’s AutoGen framework is inspired by Marvin Minsky’s “Society of Mind” theory, which posits that intelligence emerges from the interactions of simple, mindless agents working together. In AutoGen, the `Society of Mind Agent` can orchestrate and wrap `GroupChat` and `GroupChatManager` objects that host agents and debates. Externally, it functions as a singular cohesive agent and the conversation within `GroupChat` under the `SocietyOfMindAgent` can be considered an “Inner Dialogue”. This internal discourse allows for the weighing of options and the integration of multiple viewpoints, ultimately guiding behavior and thought processes.

3.2 Code

The code introduced a dynamic hierarchy system where a hierarchy string (e.g., “ABB”) defines the structure of a group of agents. Each letter in the hierarchy string represents the nodes that consists of agents such as the Prompt Generator, Counter, Working Agent, Checker, and Subordinate Nodes. The hierarchy system allows a clear definition of the agent hierarchy, a flexible representation of the topology, and an optimized information flow.

3.2.1 Agents

Prompt Generator: The Prompt Generator is responsible for structuring and reformatting the problem before any agent attempts to solve it. It does not solve the problem itself but ensures that agents receive a clear and well-defined prompt.

Counter: The Counter Agent is responsible for tracking the number of debate rounds and ensuring the conversation does not continue indefinitely. It does not contribute to the discussion but simply counts rounds and signals progress.

Working Agent: The Working Agents are responsible for solving the problem based on the structured prompt. They generate initial solutions, review responses from other agents, and refine their answers through debate.

Checker: The Checker Agent is responsible for evaluating whether the Subordinate Agents' answers are converged. It does not solve the problem—instead, it decides whether the answers agree or if further rounds are needed. It can terminate or continue the debate based on the convergence of responses.

Subordinate Nodes: The Subordinate Nodes are nested within the SocietyOfMindAgent(s). These consist of the Prompt Generator, Counter, Working Agents, Checker, and Subordinate Nodes if these nodes host further subnodes.

3.2.2 Hierarchy String Conversion Rules

The hierarchy string needs to be processed by a function to convert it into a multi-agent topology. The function is *parse_subnodes_generic(letters)* that takes a string of letters as input and returns a list of subnodes. The function follows these rules:

Rule 1: Identifying Generations (Levels of Hierarchy)

The hierarchy is built by identifying distinct levels (generations) of agents, based on letter occurrences. The earliest occurring letter (lexicographically lowest) (e.g., 'A') is designated as Generation 0. Subsequent letters define child nodes of the previous generation, forming hierarchical relationships. Each new letter appearing after a previous letter indicates a transition to a lower level (subordinate node).

For instance, the string is "ABC":

$$"ABC" \rightarrow Gen_0 : A \rightarrow Gen_1 : B(\text{Subordinate to } A) \rightarrow Gen_2 : C(\text{Subordinate to } B)$$

Rule 2: Grouping Subordinate Nodes

For each identified generation, consecutive occurrences of the same letter are grouped under a common parent node. These grouped nodes operate within the same hierarchical layer and are peers within the structure.

Rule 3: Parent-Child Relationships

A letter that appears after another letter of the same or a lower level is assigned as its child. A new letter that appears after the lowest-ranked letter (earliest in order) indicates the start of a new group.

Rule 4: Tracking Generations Dynamically

The function tracks the depth of each letter in the hierarchy dynamically, ensuring that sibling nodes are grouped, child nodes are properly assigned to the correct parent, and the function maintains state-aware traversal to identify correct nesting.

Rule 5: Handling Nested Hierarchies (Recursive Parsing)

If a letter appears at a deeper level than its immediate predecessor, it belongs to a subordinate node. The function recursively parses the hierarchy to establish nested societies of mind.

Rule 6: Assigning Subnode Identifiers

Subordinate nodes are labeled with unique identifiers based on their structure. If a node is part of a repeated pattern, a numerical suffix is assigned to distinguish between different instances. If the string is "ABCCBCC", the subnode identifiers would be "ABB1" representing the first instance of BB under A, and "BCC1", and "BCC2" indicating two separate groups of BCC under different parent nodes.

Example:

Given the hierarchy string "ABCCBCCCABCDDCBCCDEEDEEC" and "ABCDEEDCDDDBC-CBCDDCCB", the function would parse the hierarchy like these network diagrams in Figure 3 and Figure 4

3.2.3 Workflow

The system operates by:

1. **Parsing the Hierarchy String:** The hierarchy string is processed to identify generations (Gen_0, Gen_1, ..., Gen_i ..., Gen_n), where each generation corresponds to a distinct group of agents.
2. **Forming Subnodes:** Each subnode is treated as an independent GroupChat, hosting its own set of agents (Prompt Generator, Counter, Working Agents, and Checker). These agents

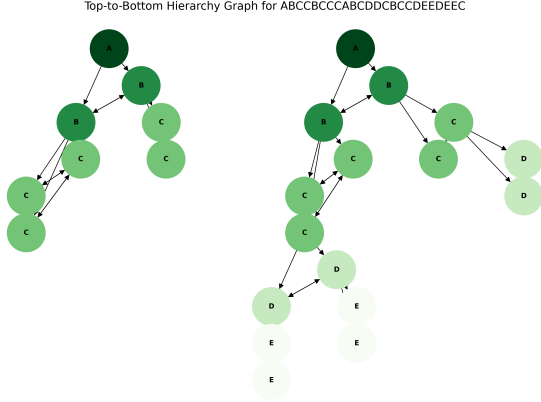


Figure 3: Visualization of the hierarchy string "ABCCBCCCABCDDCBCCDEEDEEC"

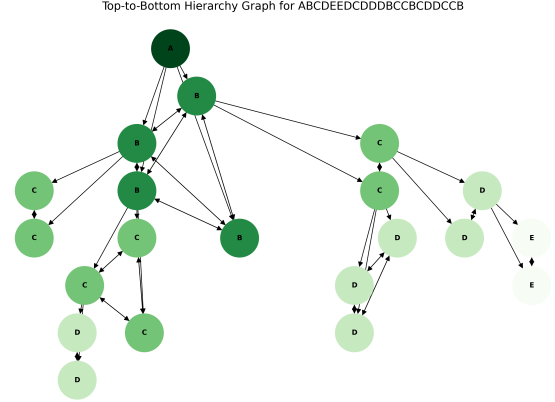


Figure 4: Visualization of the hierarchy string "ABCDEEDCDDDBCCBCDDCCB"

engage in structured discussions within their designated subnode. The nodes/subnodes are created dynamically based on the hierarchy string.

3. **Nested Debate and Refinement:** The debate starts at the lowest level (Gen_0), where Working Agents produce initial solutions. As discussions progress, higher-generation (Gen_i) subordinate subnodes generate the responses, leading to a hierarchical refinement of answers.
4. **Top-Down Control & Bottom-Up Aggregation:** The Supreme Hierarchy SocietyOf-MindAgent oversees the entire hierarchy, ensuring that responses from lower generations are aggregated and passed upwards, while high-level decisions cascade down to guide subnode debates.

This hierarchical grouping approach optimizes communication, prevents redundant discussions across unrelated agents, and ensures that information flows efficiently between different levels of the system. The combination of structured debates, iterative refinement, and dynamic agent organization enables the framework to produce well-reasoned, high-quality, accurate responses.

4 Experiment Setup and Testing Methodology

The objective of this experiment is to evaluate the performance of these multi-agent AI systems and observe its accuracy in solving a multitude of mathematical problems. These results will be compared to the performance of a single AI agent in order to contextualize these results and show the difference in accuracy, efficiency, and scalability of these different approaches.

This hierarchical debate system is important because LLMs are known to occasionally produce incorrect and inconsistent responses. The goal of this experimentation is to show whether the

hierarchical debate system and iterative refinement would lead to more accurate and reliable outputs. More specifically, this experiment aims to evaluate whether hierarchy depth further increases the accuracy at which the sample questions are answered.

4.1 Test 1: Smaller Hierarchy

The first test is designed to test a zero-shot AI agent against a multi-agent hierarchical system, following the "ABB" structure. This structure consists of one primary subnode, with two subordinate agents in the subnode. There are also support agents, those being the Prompt Generator, Counter, Checker, and Chat Manager.

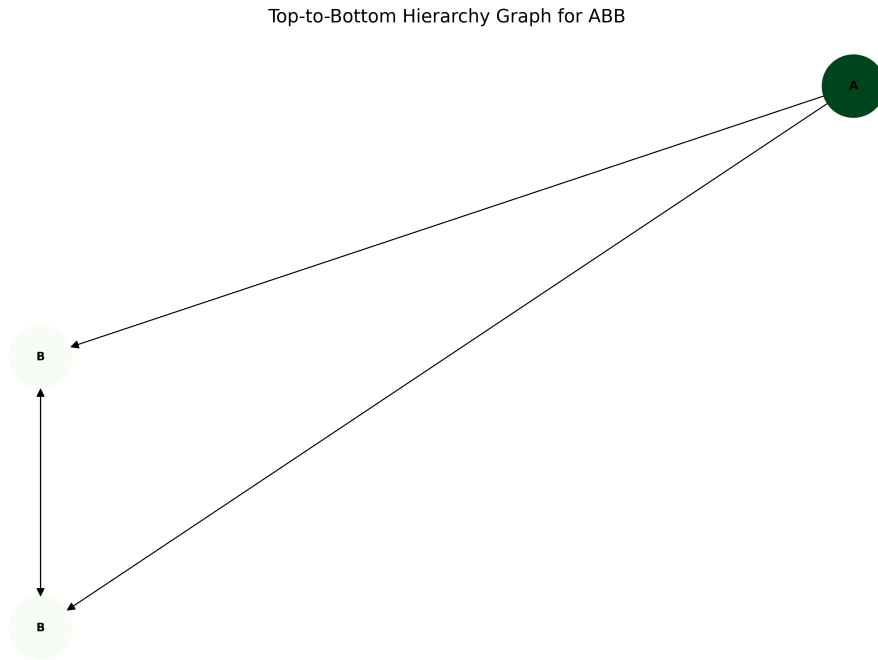


Figure 5: A Visualization of the Smaller "ABB" Hierarchy

Both the zero-shot AI agent and the "ABB" multi-agent hierarchy will be tasked with solving a set of 1000 questions, and then checked against an answer key. This experiment will establish a baseline comparison between the zero-shot agent and the simple-hierarchy multi-agent system. This will establish the effect of a small amount of collaboration between agents on solving problems.

4.2 Test 2: Large Hierarchy

Similarly to the first test, this second test will evaluate the performance of a zero-shot AI agent compared to a hierarchical system. However, this will be with the more complex hierarchical

structure "ABCCBCCC". This structure has one Gen 0 subnode, two Gen 1 subnodes, seven subordinate agents, and then the same support agents.

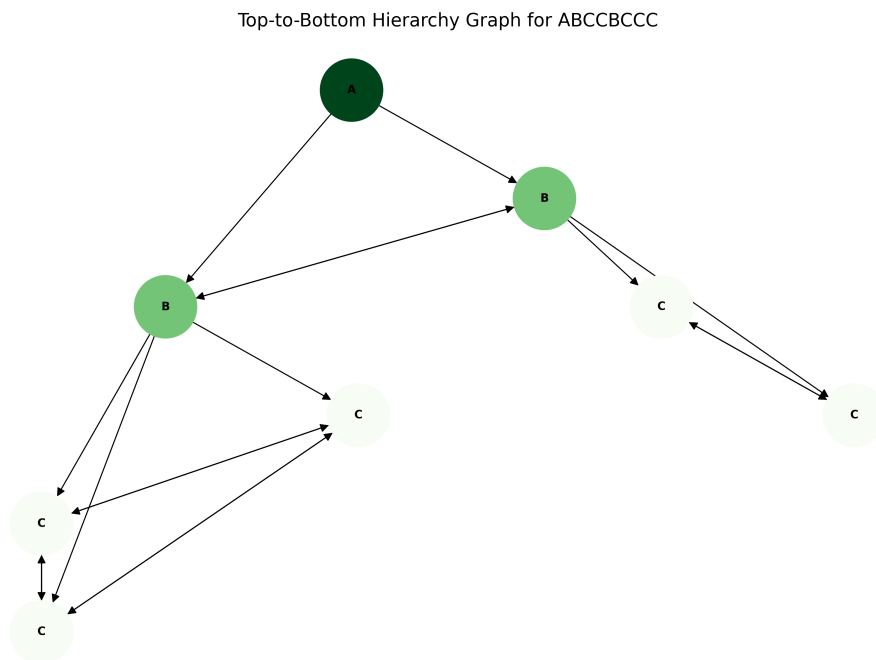


Figure 6: A Visualization of the Large "ABCCBCC" Hierarchy

In this test, both the zero-shot AI agent and the "ABCCBCCC" structure will be tasked with solving 200 questions and then checked against the answer key. While this test contains fewer questions than the first one, this is to compensate for the additional computing time due to the complex structure in this version. This test will evaluate whether the deeper hierarchies lead to better accuracy compared to the simpler hierarchies, and how that compares to the added computational intensity.

Through testing with both a simple hierarchy and a complex hierarchy, we will evaluate the practical difference in output between having a multi-agent debate structure, the complexity of the structure, and how the scalability affects the computational intensity.

4.3 Testing Procedure

The process the different AI structures followed in answering questions and providing results were done systematically, ensuring structure, consistency, and reproducibility. Each of these tests consists of question selection, execution by the agent, evaluation, and then storing the results.

4.3.1 Question Selection

The first part of the experiment is choosing a question for both structures to answer. These questions were sourced from a dataset containing mathematical problems, and also had the correct numerical answer for each question. For Test 1, we randomly selected 1000 questions, and for Test 2, we randomly selected 200 questions. Test 2 contained less questions due to the time complexity of the "ABCCBCCC" structure.

4.3.2 Agent Execution

Each question was processed in parallel by both the zero-shot agent and the multi-agent debate system. While the zero-shot agent was only prompted once to generate a response, the multi-agent model processed it in its hierarchical structure. The Prompt Generator Agent structured the problem, the Subordinate Agents independently tried to solve the problem, and the Checker Agent determined whether various answers converged. Until this convergence was achieved, the multi-agent debate continued across multiple rounds. For computational purposes and preventing an infinite loop, the maximum number of debate rounds was capped to 21.

4.3.3 Evaluation

Each answer by the AI systems were compared against the true answer provided by the dataset using an automated function. Since the problems were mathematical problems, the responses were parsed to extract only the numerical values, and this was compared against the benchmark answer. If the extracted numerical value was accurate within a tolerance of $1e-6$, then it was marked as correct. Otherwise, it was marked incorrect. The metrics that were tracked are: correct/incorrect responses, response times for both the single agent and the multi-agent system, and the round count for the multi-agent system. These are all necessary metrics to track the accuracy of both systems as well as the computational difference.

4.3.4 Data Storage

For structure and reproducibility, the test data was stored in a clear hierarchical format, ensuring easy accessibility to any necessary test data. Each test is stored in its own folder, containing: hierarchy visualization, summary of performance metrics, and stored conversation logs in another folder.

4.4 Expected Outcomes

We expect to see trade-offs between the different systems given their unique structures, across accuracy, execution time, and scalability. Accuracy can be expected to increase with the complexity of the system and the number of agents, so we would expect the "ABCCBCCC" multi-agent model

to be the most accurate, with the simpler "ABB" structure being second and the zero-shot model being the lowest. As for execution time, the reverse would be true, with the simplest zero-shot model being the fastest and the large hierarchical model being the slowest. As for scalability, we would expect there to be diminishing returns in accuracy compared to computation meaning that the difference in performance would degrade with the depth of the hierarchy.

5 Result Analysis

5.1 Comparative Analysis of Different Topological Structures

We evaluated the performance of the multi-agent system using two different hierarchical structures: *ABCCBCCC* and *ABB*. The experiments were conducted with both the multi-agent system and a zero-shot agent. Both systems utilized GPT-4o-mini-2024-07-18 and were tested on the GSM8K dataset for the task of solving mathematical word problems. The results are summarized in Table 1. The *ABCCBCCC* structure consists of eight hierarchical levels, whereas the *ABB* structure has three levels.

Table 1: Performance for Different Hierarchy Structures

Hierarchy	ABCCBCCC	ABB
Total Questions	200	1000
Accuracy (Zero-Shot Agent)	71.5%	69.8%
Accuracy (Multi-Agent System)	92.5%	90.0%
Accuracy Improvement	+21.0%	+20.2%
Time per Question (Zero-Shot)	2.73 sec	2.93 sec
Time per Question (Multi-Agent)	68.66 sec	22.48 sec
Processing Time Increase	$\times 25.15$	$\times 7.67$
Efficiency Score (Accuracy / Time)	1.35	4.00
Characteristics	Slow, High Accuracy	Balanced

The results indicate that the multi-agent system significantly enhances accuracy by 21.0% and 20.2% for the *ABCCBCCC* and *ABB* structures, respectively. Specifically, the accuracy achieved by the multi-agent system is 92.5% for the *ABCCBCCC* structure and 90.0% for the *ABB* structure, demonstrating its superiority over the zero-shot agent. Furthermore, the findings suggest that both structural complexity and the number of nodes influence accuracy, with increased complexity and a greater number of nodes leading to improved performance. However, this improvement incurs a computational cost, as the processing time increases by a factor of 25.15 for the *ABCCBCCC* structure and 7.67 for the *ABB* structure, highlighting the inherent trade-off between accuracy and computational efficiency.

To quantitatively assess this trade-off, the efficiency score, defined as the ratio of accuracy to processing time, is introduced. The ABCCBCCC structure achieves an efficiency score of 1.35, indicating a highly accurate but computationally intensive system, whereas the ABB structure attains a score of 4.00, representing a more balanced trade-off. These results suggest that while the ABCCBCCC structure prioritizes accuracy at the expense of efficiency, the ABB structure provides a more optimal balance between accuracy and processing time.

Furthermore, beyond computational efficiency, the financial cost associated with running the multi-agent system must also be considered. In this system, each node hosts multiple agents that iteratively generate responses and engage in debates. As a result, the number of requests to the LLM model increases proportionally with the number of nodes and agents, potentially leading to significantly higher costs.

5.2 Conversation Analysis

Through directly looking into the interaction between agents in the hierarchical debate structures, we can gain insights as to how the debate process contributes to the improved accuracy we observed in the experiment statistics. In addition to the debate, it is important to note the contributions of the Prompt Generator and the Checker Agent, which works to break down the problem in a defined way and ensures convergence to a single final answer. This provides structure to the debate process, ensuring consistent results.

The major advantage of the multi-agent system is its ability to self-correct through peer review, reducing the likelihood of individual errors by AI agents affecting the final results. Even if some or all AI agents may have the wrong answer initially, the debate process allows them to build on their previous attempts and collaborate, guiding them towards converging with the correct answer eventually.

5.2.1 Example 1: Quick Convergence

The following is an example of a quick convergence from the ABB hierarchy’s conversation history, testing its capability against 1000 questions. The question was the following:

"Janet's ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder at the farmers' market daily for 2 dollars per fresh duck egg. How much in dollars does she make every day at the farmers' market?"

In this example, both agents correctly calculated that there were 9 eggs left, and 18 dollars in revenue. Since they both arrived at the same conclusion, the Checker Agent recognized the convergence and terminated the debate.

5.2.2 Example 2: Delayed Convergence

While oftentimes the AI agents arrive at the same answer in the first round, there are times when they disagree, and a debate takes place until they come to an agreement. The following is an example from the BCC (part of ABCCBCCC) conversation history when tested with 200 questions. The question was:

"Josh decides to try flipping a house. He buys a house for 80,000 dollars and then puts in 50,000 dollars in repairs. This increased the value of the house by 150 percent. How much profit did he make?"

While the correct answer was 70,000 dollars, both agents initially arrived at an incorrect answer. Agent 1 incorrectly computed the total expenditure instead of the profit, declaring the answer to be 130,000. Agent 2 incorrectly added the percentage increase rather than multiplying it, declaring 200,000. The Checker Agent detected the disagreement between the two agents, and the debate continued into the next round. In the second round, Agent 1 recognized their misunderstanding by recognizing that the question specifically asked for profit and not the total cost. Meanwhile, Agent 2 also adjusted their answer after seeing Agent 1's reasoning and recognizing that they were supposed to multiply the percentage rather adding it. After this, both agents converged at the correct answer, 70,000 dollars. The Checker Agent validated that both agents came to the same answer, and terminated the debate.

This highlights the importance of the debate system and having multiple agents in arriving to an accurate answer. Even with both agents arriving at the wrong answer initially, they reviewed each other's responses and recognized their own wrongs, ultimately converging on the correct answer. With a zero-shot agent, such debate would have never taken place, and the initial wrong answers would have been considered the final. Such an example indicates that increased debate, with potentially more agents, would further help to increase the accuracy of the answers, although this would require more computational overhead.

6 Discussion

Our study suggests that incorporating hierarchical depth into the communication topology of LLM agents enhances overall accuracy by enabling multi-level collaboration and introducing redundancy to mitigate errors. While our results do not explicitly demonstrate that a deeper hierarchical communication topology (ABCCBCCC) outperforms a shallower one (ABB) in accuracy, the analysis of conversation history suggests that increased redundancy improves self-correction. This effect may become more pronounced when testing these systems on more complex problem sets.

However, redundancy and computational efficiency are inherently trade-offs, and implementing redundancy within a multi-agent debate framework is computationally demanding.

Looking ahead, we propose exploring adaptive optimization methods—such as evolutionary algorithms (e.g., NEAT[12]) or reinforcement learning—to optimize the communication topology to specific problem characteristics, although defining what constitutes a “more optimized” topology remains an open question given the high-dimensional nature of intelligence. In addition, dynamically modifying intrinsic agent properties, such as character or temperature, would promote divergent, out-of-the-box reasoning by diversifying response generation[13]; however, this approach requires regenerating outputs for evaluation within an evolutionary framework, further increasing computational demands.

7 Conclusion

This study demonstrates that incorporating a hierarchical communication topologies into a multi-agent debate system enhances accuracy compared to zero-shot reasoning. Furthermore, our analysis suggests that redundancy in communication may strengthen the system’s self-correcting ability, leading to more reliable outcomes.

References

- [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.
- [2] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018.
- [3] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI*, 2019. Accessed: 2024-11-15.
- [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- [5] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020.
- [6] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao

Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.

- [7] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners, 2023.
- [8] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback, 2023.
- [9] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023.
- [10] Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate, 2023.
- [11] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation, 2023.
- [12] Kenneth O. Stanley and Risto Miikkulainen. Evolving neural networks through augmenting topologies. *Evolutionary Computation*, 10(2):99–127, 2002.
- [13] 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, 2024.

8 Appendix: Code

```
1 from autogen import UserProxyAgent, ConversableAgent, GroupChat, GroupChatManager
2 from autogen.agentchat.contrib.society_of_mind_agent import SocietyOfMindAgent
3 import datetime
4 from dotenv import load_dotenv
5 import json
6 import matplotlib.pyplot as plt
7 import networkx as nx
8 from networkx.drawing.nx_agraph import graphviz_layout
9 import openai
10 import os
11 from pathlib import Path
12 import re
13 import sys
14 import time
15
16
17 # The name of the GPT model to be used in this multi-agent system
18 gpt_model = "gpt-4o-mini-2024-07-18"
19
20 def acquire_oak():
21     """
22     Attempts to load the OpenAI API key from a '.env' file. If the key
23     is not found, the user is prompted to enter it manually.
24
25     Returns:
26         str: The OpenAI API key to be used for subsequent calls.
27     """
28     load_dotenv() # Load environment variables from .env if present
29     api_key = os.getenv("OPENAI_API_KEY")
30     if not api_key:
31         print(".env File Not Found.")
32         # Repeatedly ask user for API key until a non-empty value is provided
33         while True:
34             try:
35                 api_key = input("Enter OpenAI API Key Here: ")
36                 if not api_key:
37                     raise ValueError("API Key cannot be empty.")
38                 break
39             except Exception as e:
40                 print(f"Error: {e}\nPlease enter the API Key again.")
41     else:
42         print("OpenAI API Key Found!")
43
44     # Set the global OpenAI API key once obtained
45     openai.api_key = api_key
46     return api_key
47
48
49 def validate_oak():
50     """
```

```

51     Validates that the acquired OpenAI API key is usable by attempting
52     to list available models from the OpenAI API. Exits the script if
53     validation fails.
54     """
55     try:
56         openai.models.list()
57         print("API Key is valid    ")
58     except Exception as e:
59         print(f"Error validating API key: {e}")
60         sys.exit(1)
61
62
63 def llm_configurator(api_key):
64     """
65     Constructs a configuration dictionary for the language model,
66     containing the model name and the provided API key.
67
68     Args:
69         api_key (str): The OpenAI API key to be used.
70
71     Returns:
72         dict: A dictionary specifying the model and API key settings.
73     """
74     return {
75         "model": gpt_model,
76         "api_key": api_key
77     }
78
79
80 def initiate_user_query():
81     """
82     Prompts the user for a problem or query.
83
84     Returns:
85         str: The user-input query string.
86     """
87     query = input("Please enter your query/problem here: ")
88     return query
89
90
91 def get_integer_input(prompt, min_value=None, error_message="Invalid input."):
92     """
93     Continuously prompts the user for an integer input until a valid value is
94     provided.
95
96     Args:
97         prompt (str): The message to display when asking for user input.
98         min_value (int, optional): The minimum valid integer value. If provided,
99             any
100             user input less than or equal to this value is
101             rejected.
102         error_message (str, optional): The message to display if validation fails.

```

```

100
101 Returns:
102     int: The valid integer provided by the user.
103     """
104 while True:
105     try:
106         value = int(input(prompt))
107         if min_value is not None and value <= min_value:
108             print(error_message)
109             continue
110         return value
111     except ValueError:
112         print(error_message)
113
114
115 def get_yes_no_input(prompt):
116     """
117     Prompts the user for a yes/no response (Y/N). Repeats until the user
118     enters a valid response.
119
120     Args:
121         prompt (str): The question or instruction to show the user.
122
123     Returns:
124         bool: True if the user inputs a "yes"-like response, False if "no"-like.
125     """
126 while True:
127     value = input(prompt).strip().lower()
128     if value in ['y', 'yes']:
129         return True
130     elif value in ['n', 'no']:
131         return False
132     else:
133         print("Please enter 'Y' or 'N'.")
134
135
136 def extract_numeric_answer(text):
137     """
138     Extracts a numeric answer (float) from a string, specifically if the string
139     contains a line matching the pattern '#### <number>'.
140
141     Args:
142         text (str): The text that potentially contains the numeric answer.
143
144     Returns:
145         float or None: The parsed float value if found, otherwise None.
146     """
147     match = re.search(r"####\s*([\d.]+)", text)
148     return float(match.group(1)) if match else None
149
150
151 def evaluate_numeric_answer(system_answer, benchmark_answer, tolerance=1e-6):

```

```

152     """
153     Evaluates whether the multi-agent system's numeric answer matches a benchmark
154     numeric answer within a given tolerance.
155
156     The expected format of both answers is: "#### <Numeric Value>"
157
158     Args:
159         system_answer (str): The final answer from the multi-agent system.
160         benchmark_answer (str): The reference/benchmark answer to check against.
161         tolerance (float, optional): Tolerance for floating-point comparison.
162             Default is 1e-6.
163
164     Returns:
165         bool: True if the numeric values match within the given tolerance; False
166             otherwise.
167     """
168     # Extract numeric answers from both system and benchmark
169     system_value = extract_numeric_answer(system_answer)
170     benchmark_value = extract_numeric_answer(benchmark_answer)
171
172     # If either extraction fails, we cannot evaluate correctly
173     if system_value is None or benchmark_value is None:
174         print("Error: Could not extract numeric answer from one or both responses.")
175         return False
176
177     # Compare the absolute difference to the specified tolerance
178     if abs(system_value - benchmark_value) < tolerance:
179         return True
180     else:
181         return False
182
183 def parse_subnodes_generic(letters):
184     """
185     Parses a sub-graph's list of letters into a hierarchical grouping (i.e.,
186     generations).
187
188     The main logic is as follows:
189     1. Identify the lowest-ranked letter (e.g., 'A' if it exists).
190     2. Gather all letters of the next rank (e.g., 'B') until the next occurrence
191        of the lowest-ranked letter.
192        This collection becomes "Gen_0" in the output.
193     3. For subsequent ranks (B, C, D, ...), gather consecutive letters of the
194        next rank and store those
195        in the appropriate "Gen_X" list.
196
197     For example, in the string "ABCCBCCC":
198     - 'A' is the lowest letter. You collect 'B' letters that appear after 'A'
199       but before the next 'A'.
200     - Then for each 'B', you gather consecutive 'C's, and so on.

```

```

197 Args:
198     letters (list or str): A sequence of letters forming a sub-graph. For
199                             instance:
200                                 ['A', 'B', 'C', 'C', 'B', 'C', 'C'] or "ABCCBCCC".
201
202 Returns:
203     dict: A dictionary containing the following keys:
204         "nodes" (list): The raw sequence of letters in the sub-graph.
205         "Gen_0", "Gen_1", "Gen_2", ... : Lists of "subnode" identifiers
206                                     extracted from the sequence.
207         Each key-value pair corresponds to letters arranged in generations.
208
209     Example return structure:
210     {
211         "nodes": ["A","B","C","C","B","C","C"],
212         "Gen_0": ["ABB1"],
213         "Gen_1": ["BCC1", "BCC2"],
214         "subnode_counts": { ... }
215     }
216
217 """
218 # Convert the input to a string if it is a list of letters
219 s = letters if isinstance(letters, str) else "".join(letters)
220 if not s:
221     # If no letters are provided, return an empty structure
222     return {"nodes": [], "Gen_0": []}
223
224 # Determine the minimum and maximum letters in the string (e.g., A, B, C, ...)
225 min_letter = min(s)
226 max_letter = max(s)
227
228 # Convert those letters to numerical "ranks" based on 'A' = 0, 'B' = 1, etc.
229 min_rank = ord(min_letter) - ord('A')
230 max_rank = ord(max_letter) - ord('A')
231
232 # Prepare a dict to store the final results
233 result = {"nodes": list(s)} # raw structure (the exact letters in the sub-
234                             graph)
235
236 # -----
237 # PART A: Handle rank 0 (lowest letter) logic
238 # -----
239 letter_r0 = chr(ord('A') + min_rank) # e.g. 'A' if min_letter is 'A'
240 letter_r1 = chr(ord('A') + min_rank + 1) # e.g. 'B' if the next rank exists
241
242 gen_0_key = f"Gen_{0}"
243 result[gen_0_key] = []
244
245 # Gather positions in the string where the rank_0 letter occurs
246 positions_of_r0 = [i for i, ch in enumerate(s) if ch == letter_r0]
247
248 if positions_of_r0:
249     # Append a sentinel index at the end of the string
250     positions_of_r0.append(len(s))

```

```

246
247 # For each occurrence of the rank_0 letter (e.g. 'A'), collect
248 # consecutive letters of rank_1 (e.g. 'B') until we hit the next 'A'
249 for idx in range(len(positions_of_r0) - 1):
250     start = positions_of_r0[idx]
251     end = positions_of_r0[idx + 1]
252
253     # Count how many letters of the next rank appear between these indices
254     B_count = 0
255     if min_rank + 1 <= max_rank:
256         B_count = sum(1 for j in range(start + 1, end) if s[j] ==
257                        letter_r1)
258
259     # Construct a subnode string (e.g. "A" + "B"*some_count)
260     subnode_str = letter_r0 + (letter_r1 * B_count)
261
262     # Tally how many times this particular subnode string has appeared
263     subnode_count = result.get("subnode_counts", {})
264     subnode_count[subnode_str] = subnode_count.get(subnode_str, 0) + 1
265     # Generate a unique identifier for this subnode
266     unique_subnode_id = f"{subnode_str}{subnode_count[subnode_str]}"
267     result[gen_0_key].append(unique_subnode_id)
268     result["subnode_counts"] = subnode_count
269
270 # For subsequent generations: B -> C, C -> D, etc.
271 for i in range(min_rank, max_rank):
272     letter_current = chr(ord('A') + i)
273     letter_next = chr(ord('A') + i + 1)
274     # Skip the rank_0 subnodes (already handled)
275     if i == min_rank:
276         continue
277
278     # Prepare a new list to store subnodes for this generation
279     result_gen_key = f"Gen_{i - min_rank}"
280     result[result_gen_key] = []
281
282     # Scan through the string to find letter_current (e.g. 'B'),
283     # then gather consecutive letter_next (e.g. 'C') immediately following
284     pos = 0
285     while pos < len(s):
286         if s[pos] == letter_current:
287             pos2 = pos + 1
288             count_next = 0
289             while pos2 < len(s) and s[pos2] == letter_next:
290                 count_next += 1
291                 pos2 += 1
292
293             if count_next > 0:
294                 # Construct e.g. "B" + "C"*count
295                 subnode_str = letter_current + (letter_next * count_next)
296                 result[result_gen_key].append(subnode_str)

```

```

297         pos = pos2
298     else:
299         pos += 1
300
301     return result
302
303
304 def generate_hierarchy_graph(hierarchy_string):
305     """
306     Creates and saves a NetworkX-directed graph for visualizing the subnode
307     hierarchy
308     derived from a given hierarchy string.
309
310     This function:
311     1. Splits the hierarchy string by top-level 'A' occurrences (if any).
312     2. For each top-level node, calls parse_subnodes_generic() to identify
313        subnodes.
314     3. Builds a Directed Graph (nx.DiGraph) where each subnode is linked
315        according to
316        a parent-child (and sibling) relationship.
317     4. Visualizes the resulting graph with matplotlib, applying a color gradient
318        based on the rank of each letter ('A' as rank=0, 'B' as rank=1, etc.).
319     5. Saves the resulting graph to "hierarchy_graph.png" in the global
320        output_directory.
321
322     Args:
323     hierarchy_string (str): A string that starts with 'A', representing
324        the structure of your multi-agent hierarchy.
325        Example: "ABB" or "ABCCBCCC".
326
327     Raises:
328     ValueError: If the hierarchy string does not start with 'A'.
329
330     Returns:
331     None. The generated graph image is saved to disk.
332     """
333     # The hierarchy string must begin with 'A'
334     if not hierarchy_string.startswith("A"):
335         raise ValueError("The hierarchy string must start with 'A'.")
336
337     # 1. Split the hierarchy string into top-level nodes
338     nodes = []
339     current_node = []
340     for letter in hierarchy_string:
341         # Each time we hit an 'A', start a new top-level node
342         if letter == "A":
343             if current_node:
344                 nodes.append(current_node)
345             current_node = []
346         current_node.append(letter)
347     if current_node:
348         nodes.append(current_node)

```

```

345
346 print("===== HIERARCHY STRUCTURE
      =====")
347 print(f"[INFO] Found {len(nodes)} top-level node(s) in the string.")
348
349 # Parse each top-level node into subnode structures
350 all_subnodes = []
351 for idx, node_letters in enumerate(nodes, start=1):
352     print(f" - Node {idx} has letters: {node_letters}")
353
354     # Parse the subnode structure of this node
355     subnode_structure = parse_subnodes_generic(node_letters)
356     all_subnodes.append(subnode_structure)
357
358     # Print out any discovered generation keys (e.g., 'Gen_0', 'Gen_1', etc.)
359     for gen_key in sorted(k for k in subnode_structure.keys() if k.startswith(
        "Gen_")):
360         gen_list = subnode_structure[gen_key]
361         if gen_list:
362             print(f" - Node {idx} has {gen_key} subnodes: {gen_list}")
363 print("
      =====\n
      n")
364
365 # -----
366 # 3. Build a Directed Graph (DiGraph) using a stack-based logic
367 # -----
368 G = nx.DiGraph()
369 node_counter = 1
370 generation_dict = {}
371
372 # For each top-level node, build sub-graphs by linking letters in a parent-
    child fashion
373 for i, node in enumerate(nodes):
374     parent_stack = []
375     instance_count = {}
376     parent_to_children = {}
377
378     print(f"=== Building Sub-Graph for Top-Level Node {i + 1} ===")
379
380     for letter in node:
381         # Track how many times we have seen this letter so far
382         instance_count.setdefault(letter, 0)
383         instance_count[letter] += 1
384         # Create a unique node_id combining letter, instance count, and
            node_counter
385         node_id = f"{letter}_{instance_count[letter]}_{node_counter}"
386         G.add_node(node_id, label=letter)
387
388         current_rank_val = ord(letter) - ord('A')
389

```



```

390         # Find the correct parent by popping from parent_stack until we find
391         the proper rank
392     while parent_stack:
393         top_node_id = parent_stack[-1]
394         top_label = G.nodes[top_node_id]['label']
395         top_rank_val = ord(top_label) - ord('A')
396         # If parent is exactly one rank above, we link them
397         if top_rank_val == current_rank_val - 1:
398             G.add_edge(top_node_id, node_id)
399             parent_to_children.setdefault(top_node_id, []).append(node_id)
400             generation_dict[node_id] = generation_dict[top_node_id] + 1
401             break
402         else:
403             parent_stack.pop()
404
405     # If no valid parent found in the stack, we set its generation to 1
406     if not parent_stack:
407         generation_dict[node_id] = 1
408
409     # Create sibling edges if multiple children share the same parent
410     if parent_stack:
411         top_node_id = parent_stack[-1]
412         siblings = parent_to_children.get(top_node_id, [])
413         for sibling in siblings:
414             if sibling != node_id:
415                 G.add_edge(node_id, sibling)
416                 G.add_edge(sibling, node_id)
417
418     # Push the current node_id onto the stack
419     parent_stack.append(node_id)
420
421     print(f"Finished building sub-graph for Node {i + 1}, node_counter={
422         node_counter}")
423     node_counter += 1
424
425     # -----
426     # 4. Visualize the Graph with color-coding by rank
427     # -----
428     # We assume 'output_directory' is a global variable or defined externally
429     # to store output images. The user can customize the color map or dot layout.
430
431     # 1) Compute the rank of each node from 'A'=0, 'B'=1, etc.
432     all_ranks = [ord(G.nodes[n]['label']) - ord('A') for n in G.nodes]
433     min_rank = min(all_ranks)
434     max_rank = max(all_ranks)
435     n_ranks = max_rank - min_rank + 1
436
437     # 2) Choose a matplotlib colormap (e.g., Greens)
438     cmap = plt.cm.Greens
439
440     # 3) Assign each node a color based on its rank, normalizing to [0..1]
441     color_list = []

```

```

440     for node_id in G.nodes():
441         letter = G.nodes[node_id]['label']
442         rank_val = ord(letter) - ord('A')
443         if n_ranks > 1:
444             normalized_val = (rank_val - min_rank) / (n_ranks - 1)
445         else:
446             normalized_val = 0
447         # invert or use directly, e.g., 1 - normalized_val if you want 'A' darkest
448         color_list.append(cmap(1 - normalized_val))
449
450     # Use graphviz's 'dot' layout for hierarchical visualization
451     plt.figure(figsize=(12, 8))
452     pos = graphviz_layout(G, prog="dot")
453
454     # Extract labels from node attributes
455     labels = nx.get_node_attributes(G, "label")
456
457     # Draw the graph with relevant parameters
458     nx.draw(
459         G, pos,
460         with_labels=True,
461         labels=labels,
462         node_size=3000,
463         node_color=color_list,
464         font_size=10,
465         font_weight="bold",
466         arrowsize=15
467     )
468
469     plt.title(f"Top-to-Bottom Hierarchy Graph for {hierarchy_string}", fontsize
470             =16)
471
472     # Save the graph to a file
473     graph_path = output_directory / "hierarchy_graph.png"
474     plt.savefig(graph_path, bbox_inches='tight', dpi=300)
475     plt.close()
476     print(f"Saved hierarchy graph to: {graph_path}\n")
477
478 def single_shot_agent(question):
479     """
480     Creates a zero-shot agent that provides an immediate answer
481     in a single response without any debate. This serves as a baseline
482     or benchmark agent for comparison with the multi-agent debate system.
483
484     Args:
485         question (str): The user-provided question or prompt to be answered.
486
487     Returns:
488         str: The raw text response from the zero-shot agent, which should
489             include a numeric answer in the format "#### <Answer>" if needed.
490     """

```

```

491 # Create an OpenAI client instance to handle the request
492 client = openai.OpenAI()
493 # Send a request to the model with system/user messages, specifying the zero-
    shot approach
494 single_shot_agent = client.chat.completions.create(
495     messages=[
496         {
497             "role": "system",
498             "content": """
499             You are a helpful assistant and a Zero-Shot Agent.\n\n
500
501             Your task is to directly answer the given question in a single
                response without any further debate or
502             interaction.\n
503             You must provide the most accurate answer possible in your first
                and only attempt.\n\n
504
505             Make sure to include '#### <Answer>' at the end of your response,
                where <Answer> is the numeric answer,
506             if the problem requires a numeric result.\n
507             There can be one, and only one, '#### <Answer>' at the end.\n
508             Do not include any commas, apostrophes, or units of measurement in
                the numerical answer.\n
509             Ensure that your final answer is actually what is being required,
                rather than an intermediate calculation!\n
510             Do not provide explanations beyond the final answer.
511             """
512         },
513         {
514             "role": "user",
515             "content": f"Question: {question}"
516         }
517     ],
518     model=gpt_model
519 )
520 # Extract the text of the response
521 single_agent_answer = single_shot_agent.choices[0].message.content
522 return single_agent_answer
523
524
525 def spawn_user_interface(llm_config):
526     """
527     Creates a user-facing agent (UserProxyAgent) that
528     initiates conversations with the hierarchy.
529
530     Args:
531         llm_config (dict): The language model configuration (e.g. model name, API
                    key).
532
533     Returns:
534         UserProxyAgent: The agent that represents the user side of the
                    conversation,

```

```

535         forwarding user queries to the system.
536     """
537     user_interface = UserProxyAgent(
538         name="User_Interface",
539         # This function checks for the word 'TERMINATE' in the content to decide
           if it should stop
540         is_termination_msg=lambda x: x.get("content", "").find("TERMINATE") >= 0,
541         human_input_mode="NEVER", # No interactive human input during the script
542         code_execution_config=False, # No code execution permission
543         llm_config=llm_config,
544     )
545     return user_interface
546
547
548 def spawn_prompt_generator(llm_config, gen_name, subnode):
549     """
550     Spawns an agent whose sole purpose is to generate a structured prompt
551     based on the user problem statement, without actually solving it.
552
553     Args:
554         llm_config (dict): Configuration dictionary for the LLM.
555         gen_name (str): The generation label (e.g., 'Gen_0').
556         subnode (str): The identifier for this subnode (e.g., 'ABB1').
557
558     Returns:
559         ConversableAgent: An agent that will produce a prompt summarizing the
560                           problem
561                           and its requirements, instructing other agents to solve
562                           it.
563     """
564     prompt_generator = ConversableAgent(
565         name="Prompt_Generator",
566         # The following system_message instructs the agent on how to create
           prompts
567         system_message=(f"""
568             You are a prompt generator designed to assist other language model
569             agents in solving problems accurately.\n
570             In this multi-agent system, you have been assigned to a particular
571             section, called a subnode.\n
572             You are a member of subnode '{subnode}', belonging to the generation
573             '{gen_name}'.\n\n
574
575             Given a user-provided problem statement, perform the following tasks
576             without attempting to solve the problem.\n\n
577
578             Here are your instructions:\n
579             0. **YOU MUST NEVER SOLVE THE PROBLEM**\n
580             1. **Identify the Category:** Determine the category or subject of the
581                problem (e.g., Mathematics, Science,
582                Language Arts).**\n
583             2. **Summarize the Problem:** Provide a summary of the problem,
584                highlighting key information and relevant

```

```

577         data for solving it.**\n
578     3. **Specify Requirements:** Outline what is needed to solve the
579         problem. The agents must be able to solve
580         the problem based on the information that you give them, so it must be
581         complete.**\n
582     4. **Tell agents that they need to solve the problem.**\n
583     5. **If the debate continues because convergence is not reached,
584         simply repeat the original prompt.**\n
585     6. You are NEVER allowed to say 'TERMINATE'.\n\n
586
587     Generate the prompt in English following this structure:\n\n
588
589     ---\n"
590     **Category:** [Determined Category]\n\n
591     **Problem Summary:** [Summary]\n\n
592     **Requirements:** [List of Requirements]\n
593     **Instructions to Agents:** [Clear instructions to solve the problem]
594     ---
595     """),
596     is_termination_msg=lambda x: x.get("content", "").find("TERMINATE") >= 0,
597     code_execution_config=False,
598     llm_config=llm_config,
599 )
600 return prompt_generator
601
602 def spawn_counter(llm_config, gen_name, subnode):
603     """
604     Spawns an agent that simply counts the number of debate rounds until '
605     TERMINATE' occurs.
606
607     Args:
608         llm_config (dict): The LLM configuration.
609         gen_name (str): Generation label (e.g., 'Gen_0', 'Gen_1').
610         subnode (str): Identifier for this subnode.
611
612     Returns:
613         ConversableAgent: An agent that will output "Round X" each time it speaks,
614         incrementing the round count each time, until
615         termination.
616
617     """
618     counter = ConversableAgent(
619         name="Counter",
620         # The system message explains the agents role in counting debate rounds
621         system_message=(
622             "You are a helpful assistant. You are one of the agents of a multi-
623             agent debate system.\n\n"
624             "In this multi-agent system, you have been assigned to a particular
625             section, called a subnode.\n"
626             f"You are a member of subnode '{subnode}', belonging to the generation
627             '{gen_name}'.\n\n"
628             "A debate is composed of a series of responses by agents.\n"

```

```

621         "A debate round begins with the response by the first agent and ends
        with the evaluation by the Checker agent.\n"
622     "You are going to count the number of debate rounds.\n"
623     "If you do not see any response yet, then it is the beginning of the
        first round, labelled Round 1.\n\n"
624     "***Your output must follow the following format: 'Round <<round_number
        >>'**\n"
625     "***Do not restart at Round 1 unless it's a new problem.**\n"
626     "***Stop counting when 'TERMINATE' is detected.**\n"
627     "***Do not answer the prompt by the user. Your sole job is counting
        rounds, following the instructions above.**"
628     ),
629     llm_config=llm_config
630 )
631 return counter
632
633
634 def spawn_subordinate_agents(llm_config, gen_name, subnode, idx_sub, agent_letter)
635 :
636     """
637     Spawns subordinate agents that attempt to solve the problem after receiving
        the prompt.
638     They can revise their answers in subsequent rounds based on peer responses.
639
        Args:
640         llm_config (dict): The LLM configuration to use for this agent.
641         gen_name (str): The generation name, e.g. 'Gen_0'.
642         subnode (str): The subnode identifier (e.g., 'ABB1').
643         idx_sub (int): Numerical index for this subordinate agent (e.g. 1, 2, etc
            .).
644         agent_letter (str): A single letter that differentiates subordinate agents
            (e.g. 'B' or 'C').
645
        Returns:
646         ConversableAgent: A subordinate agent that will produce its own solution
            in Round 1
647                             and refine it in further rounds after reviewing others'
            solutions.
648
649     """
650     agent = ConversableAgent(
651         # Construct a name like "Subordinate_Agent_B1_Gen_0_Subnode_ABB1"
652         name=f"Subordinate_Agent_{agent_letter}{idx_sub}_{gen_name}_Subnode_{
            subnode}",
653         system_message=(f"""
654             You are a helpful assistant. You are one of the agents of a multi-
            agent debate system.\n\n
655
656             In this multi-agent system, you have been assigned to a particular
            section, called a subnode.\n
657             The subnode works like a group chat. You will be chatting with the
            other agents of your kind.\n\n
658

```

```

659         You are a Subordinate member of subnode '{subnode}', belonging to the
        generation '{gen_name}'.\n\n
660
661         Given a structured prompt provided by the Prompt Generator agent,
        please do the following:\n"
662         1. In Round 1, do NOT look at or acknowledge previous Agents'
        responses. Solve the problem INDEPENDENTLY,
663         ALWAYS providing your own explanation because the user needs to
        see your reasoning;\n
664         2. IF, AND ONLY IF, you are in Round 2 or later, review the
        responses of the other Subordinate Agents and
665         provide your own explanation based on their contributions;\n
666         3. From Round 2 onwards, refine your own response until the
        Checker Agent says 'TERMINATE'.\n\n
667
668         A debate is composed of a series of responses by agents.\n
669         A debate round begins with the response by the first agent and ends
        with the evaluation by the Checker agent.\n
670         The round ends after all agents have responded.\n
671         Please NEVER count the Rounds. There is a specific Counter Agent for
        that purpose.\n\n
672
673         Make sure to include '#### <Answer>' at the end of your response,
        where <Answer> is the numeric answer,
674         if the problem requires a numeric result. Before that, make sure to
        always provide your explanation.\n
675         There can be one, and only one, '#### <Answer>' at the end.\n
676         Do not include any commas, apostrophes, or units of measurement in the
        numerical answer.\n
677         Ensure that what you put in '#### <Answer>' is actually what is being
        required, rather than an intermediate number!\n
678         """),
679         is_termination_msg=lambda x: x.get("content", "").find("TERMINATE") >= 0,
680         code_execution_config=False,
681         llm_config=llm_config,
682         description="""
683             Subordinate Agent tasked with:\n
684             1. Generating a response to find a solution to the problem
        provided by the user;\n
685             2. Reviewing other Subordinate Agents' responses only from
        Round 2 onwards;\n
686             3. Refining its own response in future rounds of debate based
        on the contributions by the other
687             Subordinate Agents, until the Checker Agent says 'TERMINATE'.
688             """)
689     )
690     return agent
691
692
693 def spawn_checker(llm_config, gen_name, subnode):
694     """

```

```

695     Spawns the Checker agent for a given subnode. This agent evaluates all
696         subordinate agents
697     responses and decides if they converge (agree on the same answer). If they do,
698         it outputs 'TERMINATE';
699     otherwise, it outputs 'CONTINUE'.
700
701     Args:
702         llm_config (dict): The LLM configuration.
703         gen_name (str): The generation name of the subnode (e.g. 'Gen_0').
704         subnode (str): The identifier for the subnode (e.g., 'ABB1').
705
706     Returns:
707         ConversableAgent: The Checker agent that can end the debate or request
708             more rounds.
709
710     """
711     checker = ConversableAgent(
712         name="Checker",
713         system_message=(f"""
714         You are a helpful assistant. You are one of the agents of a multi-agent
715         debate system.\n\n
716
717         In this multi-agent system, you have been assigned to a particular section
718             , called a subnode.\n
719         The subnode works like a group chat. However, you do not participate to
720             the discussion.\n\n
721
722         You are a member of subnode '{subnode}', belonging to the generation '{
723             gen_name}'.\n\n
724
725         You are tasked with evaluating the responses generated by multiple LLM
726             Agents within your subnode.\n
727         Those agents will be collaborating to solve the same problem.\n
728         Your job is to evaluate their responses and determine if they have reached
729             a consensus.\n
730         If the agents provide responses that are similar in meaning and consistent
731             , then you must reply 'TERMINATE' to end the debate.\n
732         Otherwise, reply 'CONTINUE' to allow another discussion round.
733         """),
734         is_termination_msg=lambda x: x.get("content", "").find("TERMINATE") >= 0,
735         code_execution_config=False,
736         llm_config=llm_config,
737     )
738     return checker
739
740 def spawn_subnode_group_chat(all_agents):
741     """
742     Creates a GroupChat instance for a given subnode, incorporating:
743     - A Chat Manager
744     - Prompt Generator
745     - Counter
746     - Subordinate Agents

```



```

737     - Checker Agent
738
739     Args:
740         all_agents (list): A list of agent objects (ConversableAgent, etc.) that
741                             will participate in the subnode chat.
742
743     Returns:
744         GroupChat: An object that orchestrates multi-turn conversations among the
745                     given agents.
746
747     """
748     # We specify the maximum number of rounds to avoid infinite loops.
749     # The speaker_selection_method ensures that each agent gets to speak in a
750     # round-robin order.
751     subnode_group_chat = GroupChat(
752         agents=all_agents, # Agents include Chat_Manager, Prompt_Generator,
753         Counter, Subordinate Agents, and Checker
754         messages=[],
755         max_round=21, # Some reasonably large limit; can be tuned or replaced
756                       # dynamically
757         speaker_selection_method="round_robin",
758     )
759     return subnode_group_chat
760
761 def spawn_chat_manager(subnode_group_chat, llm_config, gen_name, subnode,
762 chat_manager_letter, subordinates_letters):
763     """
764     Creates and configures a GroupChatManager for the specified subnode, which
765     controls
766     how the conversation flows among the agents (Prompt Generator, Counter,
767     Subordinates, Checker).
768
769     Args:
770         subnode_group_chat (GroupChat): The GroupChat object that contains the
771                                         relevant agents for this subnode.
772         llm_config (dict): The language model configuration.
773         gen_name (str): The generation name for the subnode (e.g., 'Gen_0').
774         subnode (str): The subnode identifier (e.g., 'ABB1').
775         chat_manager_letter (str): A single letter identifying the chat manager (e
776                                   .g., 'A').
777         subordinates_letters (list): Letters that identify the subordinate agents
778                                     (e.g., ['B', 'B']).
779
780     Returns:
781         GroupChatManager: An object that orchestrates the order of speaker turns
782                           and can detect termination events.
783
784     """
785     chat_manager = GroupChatManager(
786         groupchat=subnode_group_chat,
787         name=f"Chat_Manager_{chat_manager_letter}_{gen_name}_Subnode_{subnode}",
788         system_message=(f"""

```

```

776         You are managing a Group Chat where Agents must discuss and collaborate to
777             find a solution to a problem provided
778         by the user and structured into a prompt by a prompt generator.\n\n
779
780         The order of speakers in the Group Chat **must** be the following and
781             cannot be altered:\n
782         1. **Prompt Generator** (only in the very first round).\n
783         2. **Counter Agent**.\n
784         3. **All subordinate agents** (ordered alphabetically).\n
785         4. **Checker Agent** (ONLY AFTER all working agents have responded).\n
786         5. If convergence is achieved, you must also end the debate by simply
787             saying 'TERMINATE'..\n
788             Otherwise, start again from step 1 (Prompt Generator).\n\n
789
790         **You must strictly enforce this order and ensure no agent speaks out of
791             turn.**
792
793         Your responsibilities are the following:\n
794         1. Ensure that each Agent produces a response.\n
795         2. Strictly enforce the speakers order and ensure no agent speaks out of
796             turn.\n
797         3. If the Checker Agent does not detect convergence and says 'CONTINUE',
798             allow another round of debate.\n
799         In such a case, ensure that the Agents update their responses based on the
800             responses from the previous rounds
801         of debate.\n
802         4. If the Checker agents detects convergence and says 'TERMINATE',
803             terminate the debate by saying 'TERMINATE'.
804         """),
805         is_termination_msg=lambda x: x.get("content", "").find("TERMINATE") >= 0,
806         llm_config=llm_config,
807     )
808     # Print a quick info message to identify the chat manager and subordinates
809     print(
810         f" -> Manager of the GroupChat: '{chat_manager_letter}', "
811         f"Subordinates: {list(subordinates_letters)}).\n"
812     )
813     return chat_manager
814
815 def spawn_hierarchy_manager(
816     chat_manager, llm_config, gen_name, subnode, chat_manager_letter,
817     user_query
818 ):
819     """
820     Spawns a Hierarchy Manager (SocietyOfMindAgent) responsible for handing off
821         the user query
822     to the subnode-level Chat Manager. The Hierarchy Manager enforces a one-time
823         pass of the user
824     query to avoid restarting the subnode multiple times.
825
826     Args:

```

```

817         chat_manager (GroupChatManager): The GroupChatManager that orchestrates
            the subnode.
818         llm_config (dict): Configuration dict with model and API key.
819         gen_name (str): Generation label (e.g. 'Gen_0').
820         subnode (str): Subnode identifier (e.g. 'ABB1').
821         chat_manager_letter (str): Letter identifying this chat manager.
822         user_query (str): The original user query/problem statement.
823
824     Returns:
825         SocietyOfMindAgent: The hierarchy manager that can initiate a single pass
            of the subnode chat.
826
827     """
828     hierarchy_manager = SocietyOfMindAgent(
829         name=f"Hierarchy_Manager_{chat_manager_letter}_{gen_name}_Subnode_{subnode
            }",
830         chat_manager=chat_manager,
831         is_termination_msg=lambda x: x.get("content", "").find("TERMINATE") >= 0,
832         code_execution_config=False,
833         llm_config=llm_config,
834     )
835
836     # A private boolean to indicate if the subnode chat has already run
837     hierarchy_manager._subnode_done = False
838
839     def process_handoff(*args, **kwargs):
840         """
841         Initiates the subnode chat exactly once by passing the user query to the
            chat manager.
842         If run again, it would still pass the query, but in your final code you
            might
843         guard against re-initiation if you only want it to run once.
844         """
845         print(f"Hierarchy_Manager_{chat_manager_letter} passing original user
            query to subnode {subnode}.")
846         return chat_manager.initiate_chat(chat_manager, message=user_query,
            clear_history=False)
847
848     # Overwrite the default .initiate_chat method with the custom process_handoff
849     hierarchy_manager.initiate_chat = process_handoff
850
851     return hierarchy_manager
852
853 def llm_parser(aggregated_answer):
854     """
855     Uses an additional LLM-based parser to extract the final numeric answer from
856     a possibly verbose response, returning it in the standardized format '#### <
        Answer>'.
857
858     Args:
859         aggregated_answer (str): The final combined answer from a multi-agent
            debate.

```

```

860
861 Returns:
862     str: A parsed string of the form "#### <numeric_value>", with no extra
863         text.
864
865 """
866 client = openai.OpenAI()
867 parser = client.chat.completions.create(
868     messages=[
869         {
870             "role": "system",
871             "content": """
872             You are an LLM Parser. Your job is to parse the final answer from
873             other LLM agents.\n
874             You are tasked with summarizing the answer in a single number,
875             without any extra words.\n\n
876
877             Your response should just contain a number WITHOUT any commas,
878             apostrophes, or units of measurement.\n
879             Format:\n
880             #### <<Numeric Answer>>\n\
881             """,
882         },
883         {
884             "role": "user",
885             "content": f"Answer: {aggregated_answer}"
886         }
887     ],
888     model=gpt_model
889 )
890 parsed_answer = parser.choices[0].message.content
891 return parsed_answer
892
893 def create_group_chat(llm_config, gen_name, subnode):
894     """
895     Initializes and returns the subnode group chat by creating and collecting:
896     - A Prompt Generator agent,
897     - A Counter agent,
898     - Subordinate agents,
899     - A Checker agent,
900     - A Chat Manager to orchestrate the conversation.
901
902     Args:
903         llm_config (dict): The configuration for the LLM.
904         gen_name (str): The generation name, e.g. 'Gen_0'.
905         subnode (str): A string like 'ABB1' that identifies the subnode.
906
907     Returns:
908         GroupChatManager: The chat manager that coordinates all agents for this
909             subnode.
910
911     """
912     print(f"\nCreating GroupChat for Subnode: {subnode}")

```

```

907
908     # List to hold all agent instances for this subnode
909     all_agents = []
910
911     # Create the Prompt Generator (which never solves, only outlines the problem)
912     prompt_generator = spawn_prompt_generator(llm_config, gen_name, subnode)
913     all_agents.append(prompt_generator)
914
915     # Create the Counter agent
916     counter = spawn_counter(llm_config, gen_name, subnode)
917     all_agents.append(counter)
918
919     # The first letter in subnode indicates the Chat Manager (e.g. 'A')
920     # The rest are subordinate letters (e.g. 'B', 'B') or similar
921     chat_manager_letter = subnode[0]
922     subordinates_letters = [ch for ch in subnode[1:] if ch.isalpha()] # Only
        letters
923
924     # Create subordinate agents (e.g. 'B1', 'B2', etc.)
925     working_agents = []
926     for idx_sub, agent_letter in enumerate(subordinates_letters, start=1):
927         agent = spawn_subordinate_agents(llm_config, gen_name, subnode, idx_sub,
            agent_letter)
928         working_agents.append(agent)
929
930     # Add all subordinate agents to the list
931     all_agents.extend(working_agents)
932
933     # Spawn the Checker agent and add to the list
934     checker = spawn_checker(llm_config, gen_name, subnode)
935     all_agents.append(checker)
936
937     # Create the subnode group chat with all these agents
938     subnode_group_chat = spawn_subnode_group_chat(all_agents)
939
940     # Create and return the chat manager for this subnode
941     chat_manager = spawn_chat_manager(
942         subnode_group_chat, llm_config, gen_name, subnode, chat_manager_letter,
            subordinates_letters
943     )
944
945     return chat_manager
946
947
948 if __name__ == "__main__":
949     """
950     Main entry point for the multi-agent system execution. This section handles:
951
952     1. API Key retrieval and validation.
953     2. User choice of hierarchy string (default or custom).
954     3. User choice of using a JSON file with math problems or a custom query.

```

```

956 4. Setup of the multi-agent framework, including:
957     - Parsing of the hierarchy string into subnodes.
958     - Generating and saving a hierarchy graph visualization.
959     - Spawning Zero-Shot and Multi-Agent systems to answer the question(s).
960 5. Orchestrating the final output:
961     - Tracking how many questions were correctly answered by Zero-Shot vs
          Multi-Agent systems.
962     - Printing a summary of results and conversation histories.
963 """
964
965 # Step 1: Configuration and API Setup
966 print("==== Welcome to the Multi-Agent System ====\\n")
967 api_key = acquire_oak() # Retrieve API key from .env or user prompt
968 validate_oak() # Validate the acquired API key
969 llm_config = llm_configurator(api_key) # Create a dict with model and API key
          settings
970
971 # Step 2: Prompt for the Hierarchy String using get_yes_no_input
972 if get_yes_no_input("Do you want to use the default hierarchy string 'ABB'? (Y
          /N): "):
973     hierarchy_string = "ABB"
974     print(f"Using default hierarchy string: {hierarchy_string}")
975 else:
976     # Let the user enter a custom hierarchy string, ensuring it starts with 'A
          ,
977     hierarchy_string = input("Enter your custom hierarchy string (e.g.,
          ABCCBCCC): ").strip().upper()
978     while not hierarchy_string.startswith("A"):
979         print("Error: The hierarchy string must start with 'A'.")
980         hierarchy_string = input("Enter your custom hierarchy string (e.g.,
          ABCCBCCC): ").strip().upper()
981     print(f"Using custom hierarchy string: {hierarchy_string}")
982
983 # Create a timestamped output directory for logs/results
984 timestamp = datetime.datetime.now().strftime("%m-%d-%Y_%H-%M-%S")
985 output_directory = Path(f"output/{timestamp}")
986 output_directory.mkdir(parents=True, exist_ok=True)
987
988 # Step 3: Prompt whether to use the JSON file with math problems or a custom
          query
989 use_json = get_yes_no_input(
990     "Do you want to use the JSON file with math problems instead of a custom
          query? (Y/N): ")
991 if use_json:
992     # Attempt to open and parse the JSON file containing problems
993     try:
994         with open("problem_sets.jsonl", "r", encoding="utf-8") as f:
995             # Each line in the file is assumed to be a valid JSON object
996             data = [json.loads(line) for line in f]
997     except FileNotFoundError:
998         print('File not found: "problem_sets.jsonl" does not exist in the
          current directory.')
```

```

999         sys.exit(1)
1000
1001     # Extract questions and answers from the JSON lines
1002     questions = [item["question"] for item in data]
1003     answers = [item["answer"] for item in data]
1004
1005     # Ask the user how many questions to process from the loaded file
1006     num_questions = get_integer_input(
1007         "How many questions do you want the multi-agent framework to solve? ",
1008         min_value=0, error_message="The number of questions must be more than
1009         0."
1010     )
1011     questions = questions[:num_questions]
1012     answers = answers[:num_questions]
1013 else:
1014     # Custom user query mode if not using the JSON file
1015     user_query = initiate_user_query() # Prompt user for a custom question
1016     questions = [user_query] # Wrap it into a list for uniform processing
1017     answers = [None] # No benchmark answer in custom mode
1018
1019 # Step 4: Parse the Hierarchy
1020 print(f"\nParsing the Hierarchy String '{hierarchy_string}'")
1021 subnode_structure = parse_subnodes_generic(hierarchy_string)
1022
1023 # Step 5: Visualize the hierarchy (saves the output graph to disk)
1024 print("Generating the hierarchy graph for visualization...\n")
1025 generate_hierarchy_graph(hierarchy_string)
1026
1027 # Step 6: Identify and sort generations in descending order (e.g. Gen_1, Gen_0)
1028 gens = sorted([k for k in subnode_structure if k.startswith("Gen_")], reverse=
1029 True)
1030 print(f"Generations found (highest first): {gens}")
1031
1032 # Step 7: Initialize data structures for the multi-agent system
1033 final_subnode_answers = {} # e.g. { "BCC": "...", "ABB1": "...", ... }
1034 manager_to_subnodes = {} # e.g. { "A": ["ABB1"], "B": ["BCC", "BCCC"], ... }
1035
1036 # The user-facing agent who interacts with the top-level aggregator (
1037 Supreme_Hierarchy)
1038 user_interface = spawn_user_interface(llm_config)
1039 chat_managers = {} # Will hold chat managers for each subnode
1040 hierarchy_managers = {} # Will hold hierarchy managers for each subnode
1041
1042 # If you want to limit the top-level rounds (like if you have 3 subnodes total
1043 ):
1044 number_of_subnodes = 1
1045
1046 # Initialize containers for Zero-Shot vs Multi-Agent answers and evaluation
1047 parsed_single_agent_answers_list = []
1048 parsed_multi_agent_answers_list = []

```

```

1046 correct_count_single_agent = 0
1047 incorrect_count_single_agent = 0
1048
1049 correct_count_multi_agent = 0
1050 incorrect_count_multi_agent = 0
1051
1052 conversation_history_single_agent = []
1053 conversation_history_multi_agent = {}
1054
1055 single_shot_durations = []
1056 multi_agent_durations = []
1057
1058 # Outer loop: iterate over each question in the list
1059 for k in range(len(questions)):
1060     print("\n" + "=" * 80)
1061     print(f"QUESTION {k + 1}:")
1062     print("=" * 80)
1063
1064     question = questions[k]
1065     benchmark_answer = answers[k] if answers[k] is not None else None
1066
1067     # ----- Zero-Shot Agent Flow
1068     -----
1069     start_single_shot = time.time()
1070
1071     # Obtain zero-shot answer and parse out the numeric portion
1072     single_agent_answer = single_shot_agent(question)
1073     parsed_single_agent_answer = llm_parser(single_agent_answer)
1074     parsed_single_agent_answers_list.append(parsed_single_agent_answer)
1075
1076     end_single_shot = time.time()
1077     single_shot_duration = end_single_shot - start_single_shot
1078     single_shot_durations.append(single_shot_duration)
1079
1080     # Print zero-shot results
1081     print(f"Zero-Shot Agent Answer to question {k + 1}: {single_agent_answer}\n")
1082     print(f"Parsed Zero-Shot Answer: {parsed_single_agent_answer}")
1083     print(f"Zero-Shot Agent Time: {single_shot_duration:.4f} seconds")
1084
1085     # Log zero-shot conversation
1086     conversation_history_single_agent.append({
1087         "Question": question,
1088         "True Solution": f"#### {answers[k]}" if answers[k] is not None else "N/A",
1089         "Zero-Shot Agent Response": single_agent_answer,
1090         "Parsed Zero-Shot Answer": parsed_single_agent_answer,
1091         "Zero-Shot Duration (s)": single_shot_duration
1092     })
1093
1094     # Evaluate zero-shot answer if we have a benchmark
1095     if benchmark_answer is not None:

```



```

1095         if evaluate_numeric_answer(parsed_single_agent_answer,
1096                                     benchmark_answer):
1097             print("Zero-Shot Agent Evaluation:      The answer matches the
1098                   benchmark answer. Good job!")
1099             correct_count_single_agent += 1
1100         else:
1101             print("Zero-Shot Agent Evaluation:      The answer does NOT match
1102                   the benchmark answer.")
1103             incorrect_count_single_agent += 1
1104
1105     # ----- Multi-Agent Flow -----
1106     start_multi_agent = time.time()
1107
1108     # For each generation (in descending order), create subnode chat managers
1109     # and hierarchy managers
1110     for gen_name in gens:
1111         subnodes = subnode_structure[gen_name]
1112         if not subnodes:
1113             continue
1114
1115         print(f"\n--- Processing {gen_name} subnodes: {subnodes}")
1116         for subnode in subnodes:
1117             # Create a manager for this subnode
1118             chat_managers[subnode] = create_group_chat(llm_config, gen_name,
1119                                                         subnode)
1120
1121             # Map the chat manager letter to this subnode (for reference if
1122             # needed)
1123             chat_manager_letter = subnode[0]
1124             if chat_manager_letter not in manager_to_subnodes:
1125                 manager_to_subnodes[chat_manager_letter] = []
1126             manager_to_subnodes[chat_manager_letter].append(subnode)
1127
1128             # Create a hierarchy manager that feeds user_query into the
1129             # subnode chat
1130             hierarchy_managers[subnode] = spawn_hierarchy_manager(
1131                 chat_manager=chat_managers[subnode],
1132                 llm_config=llm_config,
1133                 gen_name=gen_name,
1134                 subnode=subnode,
1135                 chat_manager_letter=subnode[0],
1136                 user_query=question
1137             )
1138
1139             # Initialize conversation logs for that subnode if not already
1140             # present
1141             if subnode not in conversation_history_multi_agent:
1142                 conversation_history_multi_agent[subnode] = []
1143
1144     # Step 8: Create the Supreme Society of Minds that includes the user
1145     # interface + all hierarchy managers
1146     hierarchy_managers_values = hierarchy_managers.values()

```

```

1138 hierarchy_managers_list = list(hierarchy_managers_values)
1139 hierarchy_managers_list.append(user_interface)
1140
1141 # Build the top-level GroupChat with all the hierarchy managers + user
      interface
1142 hierarchies_group_chat = GroupChat(
1143     agents=hierarchy_managers_list,
1144     messages=[],
1145     max_round=number_of_subnodes + 1, # e.g. if you have 3 subnodes
1146     speaker_selection_method="round_robin",
1147 )
1148
1149 # A manager to orchestrate the top-level group chat, deciding if/when to
      terminate
1150 hierarchies_group_chat_manager = GroupChatManager(
1151     groupchat=hierarchies_group_chat,
1152     name="Hierarchies_Group_Chat_Manager",
1153     system_message=""
1154     You have been assigned the following tasks:
1155     Step 1. Ensure that all agents generate responses.
1156     Step 2. Evaluate whether the responses converge in similarity or not.
1157     Step 3. If they do:
1158         - Aggregate the responses into a single one,
1159         - Send it to the User_Interface,
1160         - **Add 'TERMINATE'** at the end to end the entire conversation.
1161     If they do not converge, start again from Step 1.
1162
1163     You summarize the answer in a simple single sentence ...
1164     "",
1165     is_termination_msg=lambda x: x.get("content", "").find("TERMINATE") >=
        0,
1166     llm_config=llm_config,
1167 )
1168
1169 # Create the supreme aggregator (SocietyOfMindAgent) for top-level
      coordination
1170 supreme_hierarchy = SocietyOfMindAgent(
1171     chat_manager=hierarchies_group_chat_manager,
1172     name="Supreme_Hierarchy",
1173     code_execution_config=False,
1174     llm_config=llm_config,
1175 )
1176
1177 # The user interface initiates the chat with the Supreme Hierarchy using
      the question
1178 chat_result = user_interface.initiate_chat(
1179     supreme_hierarchy,
1180     message=questions[k],
1181 )
1182
1183 # Retrieve the final aggregated multi-agent answer from the Supreme
      Hierarchy

```

```

1184     multi_agent_answer = chat_result.chat_history[0]["content"]
1185     parsed_multi_agent_answer = llm_parser(multi_agent_answer)
1186     parsed_multi_agent_answers_list.append(parsed_multi_agent_answer)
1187
1188     end_multi_agent = time.time()
1189     multi_agent_duration = end_multi_agent - start_multi_agent
1190     multi_agent_durations.append(multi_agent_duration)
1191
1192     # Print multi-agent results
1193     print(f"Multi-Agent Answer: {multi_agent_answer}\n")
1194     print(f"Parsed Multi-Agent Answer: {parsed_multi_agent_answer}")
1195     print(f"Multi-Agent System Time: {multi_agent_duration:.4f} seconds")
1196
1197     # Evaluate multi-agent answer if we have a benchmark
1198     if benchmark_answer is not None:
1199         if evaluate_numeric_answer(parsed_multi_agent_answer, benchmark_answer):
1200             print("Multi-Agent Evaluation:      The system's answer matches the
1201                   benchmark answer. Good job!")
1202             correct_count_multi_agent += 1
1203         else:
1204             print("Multi-Agent Evaluation:      The system's answer does NOT
1205                   match the benchmark answer.")
1206             incorrect_count_multi_agent += 1
1207
1208     # Save the conversation history for each subnode
1209     for subnode, chat_manager in chat_managers.items():
1210         chat_history = chat_manager.groupchat.messages if hasattr(chat_manager
1211             , "groupchat") else []
1212         conversation_history_multi_agent[subnode].append({
1213             "Question": question,
1214             "True Solution": f"#### {answers[k]}" if answers[k] is not None
1215                 else "N/A",
1216             "Chat History": chat_history
1217         })
1218
1219     # ----- Final Summary -----
1220     total_questions = len(questions)
1221     print("\n" + "=" * 80)
1222     print("FINAL BENCHMARK SUMMARY")
1223     print("=" * 80)
1224     print(f"Total Questions: {total_questions}")
1225     print(f"Correct Answers (Zero-Shot Agent)      : {correct_count_single_agent}")
1226     print(f"Incorrect Answers (Zero-Shot Agent)    : {incorrect_count_single_agent}")
1227
1228     print("." * 60)
1229     print(f"Correct Answers (Multi-Agent System)     : {correct_count_multi_agent}")
1230     print(f"Incorrect Answers (Multi-Agent System)   : {incorrect_count_multi_agent}")
1231
1232     print("." * 60)
1233     if total_questions > 0:

```

```

1228     accuracy_single_agent = (correct_count_single_agent / total_questions) *
1229         100
1229     accuracy_multi_agent = (correct_count_multi_agent / total_questions) * 100
1230     print(f"Overall Accuracy (Zero-Shot Agent)    : {accuracy_single_agent:.2f
1230         }%")
1231     print(f"Overall Accuracy (Multi-Agent System)  : {accuracy_multi_agent:.2f
1231         }%")
1232     print("-" * 80)
1233     print(f"Correct Answers      (Zero-Shot Agent)    : {
1233         correct_count_single_agent}/{total_questions}")
1234     print(f"Incorrect Answers    (Zero-Shot Agent)    : {
1234         incorrect_count_single_agent}/{total_questions}")
1235     print("." * 60)
1236     print(f"Correct Answers      (Multi-Agent System)  : {
1236         correct_count_multi_agent}/{total_questions}")
1237     print(f"Incorrect Answers    (Multi-Agent System)  : {
1237         incorrect_count_multi_agent}/{total_questions}")
1238     print("-" * 80)
1239
1240     # Here we handle saving the entire conversation history of both the Zero-Shot
1240     # agent
1241     # and the Multi-Agent Debate System, as well as generating a final JSON
1241     # summary of results.
1242
1243     # Save the zero-shot conversation history to a JSON file
1244     single_shot_history_file = output_directory / "
1244         single_shot_conversation_history.json"
1245     try:
1246         with open(single_shot_history_file, "w", encoding="utf-8") as f:
1246             # Use json.dump to write the conversation history list to disk
1247             json.dump(conversation_history_single_agent, f, ensure_ascii=False,
1247                 indent=4)
1248             print(f"Saved zero-shot conversation history at {single_shot_history_file}
1248                 ")
1249     except Exception as e:
1250         print(f"Error saving zero-shot conversation history: {e}")
1251
1252     # For each subnode, save its Multi-Agent conversation history to a separate
1253     # JSON file
1254     for subnode, chat_manager in chat_managers.items():
1255         conversation_file = output_directory / f"{subnode}_conversation_history.
1255             json"
1256
1257         try:
1258             with open(conversation_file, "w", encoding="utf-8") as f:
1258                 # conversation_history_multi_agent[subnode] holds a list of dicts
1259                 # describing the conversation
1259                 json.dump(conversation_history_multi_agent[subnode], f,
1260                     ensure_ascii=False, indent=4)
1260
1261         print(f"Saved conversation history for {subnode} at {conversation_file
1261             }")

```

```

1263
1264     except Exception as e:
1265         print(f"Error saving conversation history for {subnode}: {e}")
1266
1267 # Prepare a summary dictionary that captures the key metrics and results of
1268 # the benchmark
1269 results_summary = {
1270     "total_questions": len(questions),
1271     "correct_count_single_agent": correct_count_single_agent,
1272     "incorrect_count_single_agent": incorrect_count_single_agent,
1273     "correct_count_multi_agent": correct_count_multi_agent,
1274     "incorrect_count_multi_agent": incorrect_count_multi_agent,
1275     "overall_accuracy_single_agent": round(
1276         (correct_count_single_agent / len(questions) * 100), 2
1277     ) if total_questions > 0 else 0.0,
1278     "overall_accuracy_multi_agent": round(
1279         (correct_count_multi_agent / len(questions) * 100), 2
1280     ) if total_questions > 0 else 0.0,
1281     "timestamp": timestamp,
1282     "timing_single_shot_seconds": single_shot_durations,
1283     "timing_multi_agent_seconds": multi_agent_durations,
1284     "results": []
1285 }
1286
1287 # Populate the 'results' list with per-question information
1288 for k in range(len(questions)):
1289     question = questions[k]
1290     true_solution = f"#### {answers[k]}" if answers[k] is not None else "N/A"
1291     single_shot_result = parsed_single_agent_answers_list[k]
1292     multi_agent_result = parsed_multi_agent_answers_list[k]
1293
1294     # Construct a dictionary entry summarizing the outcomes for this question
1295     results_summary["results"].append({
1296         f"Question {k + 1}": question,
1297         "True Solution": true_solution,
1298         "Answer by The Single (Zero-Shot) Agent": single_shot_result,
1299         "Answer by The Multi-Agent Debate System": multi_agent_result,
1300         "Zero-Shot Duration (s)": single_shot_durations[k],
1301         "Multi-Agent Duration (s)": multi_agent_durations[k],
1302     })
1303
1304 # Write the benchmark summary to a JSON file
1305 summary_file = output_directory / "benchmark_summary.json"
1306 try:
1307     with open(summary_file, "w", encoding="utf-8") as f:
1308         json.dump(results_summary, f, ensure_ascii=False, indent=4)
1309
1310     print(f"Benchmark summary saved at {summary_file}")
1311
1312 except Exception as e:
1313     print(f"Error saving benchmark summary: {e}")

```