
Why Do Multi-Agent LLM Systems Fail?

Anonymous Author(s)

Affiliation

Address

email

Abstract

Despite enthusiasm for Multi-Agent LLM Systems (MAS), their performance gains on popular benchmarks are often minimal. This gap highlights a critical need for a principled understanding of why MAS fail. We create MAD, the first **Multi-Agent System Failure Dataset**, to outline the failure dynamics in MAS for guiding the development of better future systems. To enable systematic classification of failures for MAD, we build the first **Multi-Agent System Failure Taxonomy (MAST)**. We develop MAST through rigorous analysis of 150 traces, guided closely by expert human annotators and validated by high inter-annotator agreement ($\kappa = 0.88$). This process identifies 14 unique modes, clustered into 3 categories: (i) specification issues, (ii) inter-agent misalignment, and (iii) task verification. To enable scalable annotation, we develop an LLM-as-a-Judge pipeline with high agreement with human annotations. We publicly release our comprehensive dataset (MAD), the MAST, and our LLM annotator to facilitate widespread research and development in MAS.

“Happy families are all alike; each unhappy family is unhappy in its own way.” (Tolstoy [1])

“Successful systems all work alike; each failing system has its own problems.” (Anonymous Authors ‘25)

1 Introduction

Recently, Large Language Model (LLM) based agentic systems have gained significant attention in the AI community [2–4] and multi-agent systems are increasingly explored in various domains [5–11]. Despite the increasing adoption of MAS, their performance gains often remain minimal compared to single-agent frameworks [12] or simple baselines like best-of-N sampling [13]. Our empirical analysis reveals 41% to 86.7% failure rate on 7 state-of-the-art (SOTA) open-source MAS detailed in Figure 5 (Appendix C). This motivates the fundamental question we address: *Why do MAS fail?* The contributions of this paper are as follows:

- We introduce and **open-source** MAD, the first large-scale MAS failure dataset with consistent annotations from 7 MAS and two model families. And MAD-human, a detailed inter-annotator study results with human labels. Together serve to facilitate research into MAS failures.
- We introduce MAST, the first empirically grounded **taxonomy of MAS failures**, providing a structured framework for defining, understanding and annotating failures.
- We develop a scalable LLM-as-a-judge **annotation pipeline** integrated with MAST for efficiently annotating MAD and enabling analysis of MAS performance, diagnosis of failure modes, and understanding of failure breakdowns.
- We demonstrate through **case studies** that failures identified by MAST often stem from system design issues, not just LLM limitations or simple prompt following, and require more than superficial fixes, thereby highlighting the need for structural MAS redesigns.

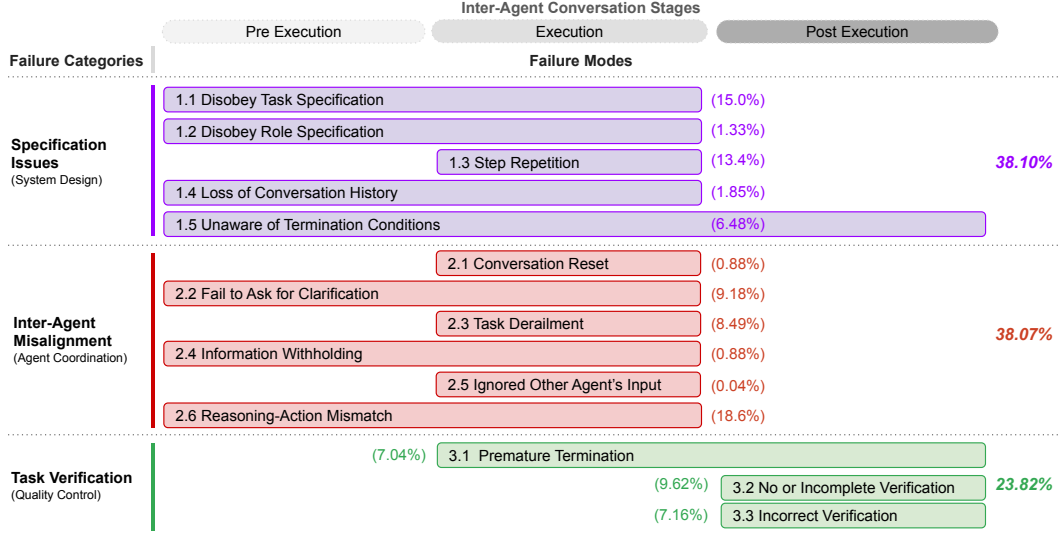


Figure 1: MAST: A Taxonomy of MAS Failure Modes. The inter-agent conversation stages indicate when a failure typically occurs within the end-to-end MAS execution pipeline. A failure mode spanning multiple stages signifies that the underlying issue can manifest or have implications across these different phases of operation. The percentages shown represent the prevalence of each failure mode and category as observed in our analysis of 1242 MAS execution traces. Detailed definitions for each failure mode and illustrative examples are available in Appendix B.

2 Related Work

The promising capabilities of agentic systems have inspired research into solving specific challenges [14–18]. While these works meaningfully contribute towards understanding specific issues or providing high-level overviews, they do not offer a fine-grained, empirically grounded taxonomy of why MAS fail across diverse systems and tasks. Numerous benchmarks also exist to evaluate agentic systems [19–24]. These evaluations are crucial but primarily facilitate a top-down perspective, focusing on aggregate performance or high-level objectives like trustworthiness and security [25, 26]. Our work complements these efforts by providing a bottom-up analysis focused on identifying specific failure modes in MAS. Several works highlight challenges in building robust agentic systems and suggest design principles, often focused on single-agent settings [13, 27]. Our work extends these insights to the multi-agent context. Despite the growing interest in LLM agents, dedicated research systematically characterizing their failure modes remains limited, particularly for MAS. While Bansal et al. [28] catalogs challenges in human-agent interaction, our contribution focuses specifically on failures within autonomous MAS execution. Other related work includes taxonomies for evaluating multi-turn LLM conversations [29] or specific capabilities like code generation [30]. These differ significantly from our goal of developing a generalizable failure taxonomy for multi-agent interactions and coordination. Further related efforts aim to improve MAS through different approaches: AgentEval [31], AGDebugger [32], MAS debugger Zhang et al. [33]. Thus, MAD and MAST represent, to our knowledge, the first empirically derived, comprehensive dataset and taxonomy focused specifically on MAS failures.

3 The Multi-Agent Systems Dataset

We introduce MAD, the Multi-Agent System Failure Dataset, an empirically grounded dataset comprising 1242 annotated execution traces collected from 7 popular MAS frameworks, covering domains of coding, math, and generic tasks. In this section, we detail our approach, which centers on building the first empirical MAS failure taxonomy, MAST, and a scalable annotation pipeline for systematic and comprehensive data collection. Figure 2 summarizes our methodological workflow.

We first collect 150 traces from five MAS frameworks, which are closely examined by six human experts. Our goal at this stage is to identify as many distinct failure modes as possible, ensuring these

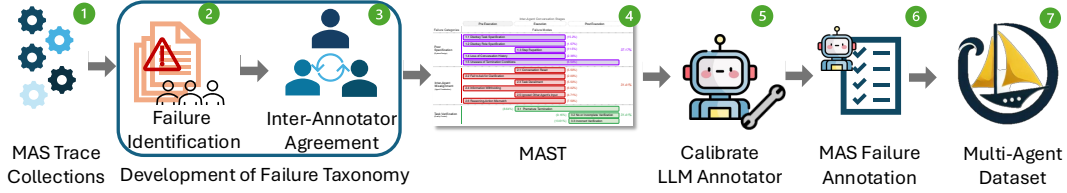


Figure 2: Methodological workflow for constructing the MAD dataset, involving the empirical identification of failure modes, the development of MAST, iterative refinement through inter-annotator agreement studies ($\kappa = 0.88$), and the creation of a scalable LLM annotation pipeline. This figure highlights our systematic approach to creating a comprehensive dataset for studying MAS failures.

observed patterns are not merely artifacts of a single system but can likely apply more broadly. To achieve this without predefined hypotheses, we adopt the **Grounded Theory** (GT) approach [34]. Five MAS frameworks and two task categories are analyzed with core GT techniques: *open coding* [35], *constant comparative analysis*, *memoing*, and *theorizing*. This iterative analysis continues until we reach *theoretical saturation*, where further data analysis does not yield new failure mode insights.

To develop a taxonomy (Section 4) that is unambiguous and consistently applicable by different annotators, we rigorously validate and refine MAST definitions through Inter-Annotator Agreement (IAA) studies. This iterative process begins with a preliminary version of MAST derived from our GT findings. In each round of IAA, three expert annotators independently label a subset of five randomly selected traces from our initial 150+ trace collection using MAST. We then facilitate discussions to collectively resolve any disagreements. Based on these discussions, we iteratively refine MAST. We measure agreement using Cohen’s Kappa score, achieving a strong average of $\kappa = 0.88$ in the final rounds.

To enable scalable and automated failure annotation for MAD, we develop an LLM-as-a-Judge pipeline (LLM annotator), building upon our validated MAST and few-shot examples from our human-annotated data (details in Appendix K), to classify observed failure modes. The LLM annotator achieves high agreement with human experts (accuracy 94%, Cohen’s Kappa of 0.77; Table 2). In order to test the generalizability of our finalized MAST and the LLM annotator, we evaluate their performance on two new MAS (OpenManus and Magentic-One) with two new benchmarks (MMLU and GAIA) not part of the initial MAST development. An additional human IAA round on these out-of-domain traces using the finalized MAST yields a strong Cohen’s Kappa score of 0.79. This demonstrates the robustness of MAST and LLM annotator. We further detail the uniqueness of MAST failure modes via a correlation study in Appendix E. Finally, we expand data collection to construct MAD, comprising 1242 annotated traces from seven popular MAS frameworks (Table 1). We also release MAD-human, consisting of all traces annotated by human experts during our IAA studies, where each annotation specifies MAST failure modes with textual justifications. We open-source MAD and MAD-human as resources to analyze MAS failure dynamics and guide robust system design.

4 The Multi-Agent System Failure Taxonomy

MAST, illustrated in Figure 1, identifies 14 fine-grained failure modes, which we map to MAS execution stages (Pre-Execution, Execution, and Post-Execution) where their root causes commonly emerge. These modes are organized into 3 overarching categories (specification issues, inter-agent misalignment, task verification) reflecting the fundamental nature of the observed failures. We do not claim MAST is exhaustive, it offers precise definitions for a structured approach to understanding why MAS fail. Detailed definitions and discussions for each failure mode (FM) are available in Appendix B, with specific examples in Appendix K.

5 Towards better Multi-Agent LLM Systems

Our analysis of MAD reveals that failure distributions differ markedly across various MAS, often reflecting their unique architectural characteristics and design philosophies. For example, as illustrated in Figure 3, we observe specific patterns: AppWorld frequently suffers from premature

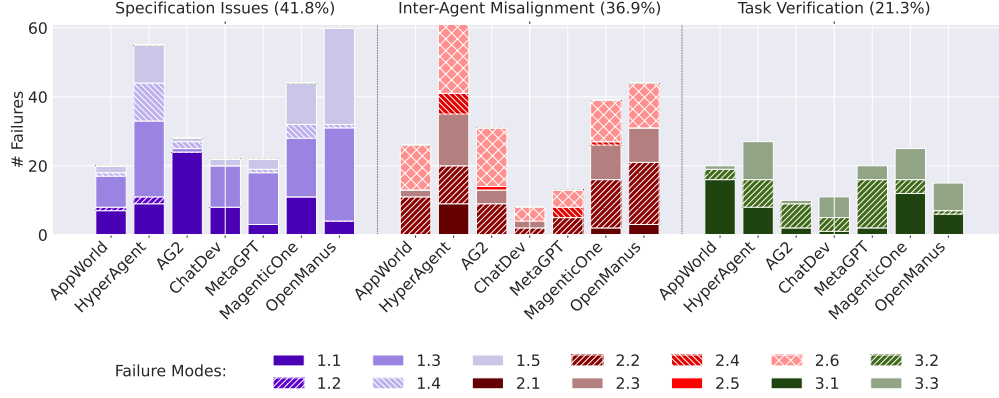


Figure 3: Distribution of failure in MAD with MAST labels on total 210 traces. This plot visualizes the failure distributions of the first 30 traces for each system. As the specific tasks and benchmarks may differ across the MAS configurations shown, these results are intended to illustrate system-specific failure profiles rather than to serve as a performance comparison across MAS.

terminations (FM-3.1), potentially due to its star topology and lack of a predefined workflow making termination conditions less obvious; OpenManus exhibits a tendency towards step repetition (FM-1.3); and HyperAgent could benefit from addressing its dominant failure modes of step repetition (FM-1.3) and incorrect verification (FM-3.3). These system-specific profiles underscore that there is no one-size-fits-all solution to MAS failures. These comparative analyses, detailed further in Appendix F, provide insights into how model choice and architectural patterns influence system performance and distribution of failures. Using our LLM annotator with MAST, developers can obtain quantitative analyses of failure profiles for specific systems. We demonstrate MAST’s practical usage in guiding MAS improvement in our case studies (Appendix H). The Failure Mode breakdown analysis (Appendix H.3) shows which failure modes were mitigated and reveals any resulting trade-offs. This granular view, moving beyond aggregate metrics, is crucial for understanding *why* an intervention works and for iterating effectively towards more robust systems.

While one could simply attribute failures in MAD to limitations of present-day LLM (e.g., hallucinations, misalignment), we conjecture that improvements in the base model capabilities will be insufficient to address the full MAST. Instead, we argue that good MAS design requires organizational understanding – even organizations of sophisticated individuals can fail catastrophically [36] if the organization structure is flawed. Previous research in high-reliability organizations has shown that well-defined design principles can prevent such failures [37, 38]. In our intervention case studies (Appendix H), we apply MAS system workflow and prompt changes respectively (results in Table 5). With the same underlying model, we achieve max improvements of 15.6%. This highlights that MAS failures can be address with better system designs.

6 Conclusion

In this study, we conduct the first systematic investigation into why MAS fail. This investigation results in the **Multi-Agent Dataset (MAD)**: a comprehensive public resource of over 1000 annotated execution traces from 7 popular MAS frameworks, which we create to outline MAS failure dynamics and guide future system development. To enable MAD’s systematic annotation and analysis, we first develop the **Multi-Agent System Failure Taxonomy (MAST)**. We build MAST through a rigorous Grounded Theory-based analysis of an initial 150 traces, validating its definitions with strong inter-annotator agreement and identifying 14 distinct failure modes across 3 categories. For scalable annotation of MAD using MAST, we then develop an LLM annotator, confirming its high agreement with human experts. Together, MAD and MAST provide a foundational framework and empirical grounding for future MAS research. We are excited about the potential of MAS, but their widespread adoption hinges on achieving greater reliability. Our work, through the public release of MAD, MAST, and the LLM annotator, contributes towards this goal. By systematically identifying and categorizing challenges, we aim to open up concrete research directions and equip the community to develop more robust and effective multi-agent systems.

References

- [1] Leo Tolstoy. *Anna Karenina*. The Russian Messenger, 1878. 139 140
- [2] Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model connected with massive apis, 2023. URL <https://arxiv.org/abs/2305.15334>. 141 142
- [3] Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. Memgpt: Towards llms as operating systems, 2024. URL <https://arxiv.org/abs/2310.08560>. 143 144 145
- [4] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Jirong Wen. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6), March 2024. ISSN 2095-2236. doi: 10.1007/s11704-024-40231-1. URL <http://dx.doi.org/10.1007/s11704-024-40231-1>. 146 147 148 149 150
- [5] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. Chatdev: Communicative agents for software development. *arXiv preprint arXiv:2307.07924*, 2023. URL <https://arxiv.org/abs/2307.07924>. 151 152 153 154
- [6] Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert Brennan, Hao Peng, Heng Ji, and Graham Neubig. Openhands: An open platform for ai software developers as generalist agents, 2024. URL <https://arxiv.org/abs/2407.16741>. 155 156 157 158 159
- [7] Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, Khaled Saab, Dan Popovici, Jacob Blum, Fan Zhang, Katherine Chou, Avinatan Hassidim, Burak Gokturk, Amin Vahdat, Pushmeet Kohli, Yossi Matias, Andrew Carroll, Kavita Kulkarni, Nenad Tomasev, Yuan Guan, Vikram Dhillon, Eeshit Dhaval Vaishnav, Byron Lee, Tiago R D Costa, José R Penadés, Gary Peltz, Yunhan Xu, Annalisa Pawlosky, Alan Karthikesalingam, and Vivek Natarajan. Towards an ai co-scientist, 2025. URL <https://arxiv.org/abs/2502.18864>. 160 161 162 163 164 165 166
- [8] Kyle Swanson, Wesley Wu, Nash L. Bulaong, John E. Pak, and James Zou. The virtual lab: Ai agents design new sars-cov-2 nanobodies with experimental validation. *bioRxiv*, 2024. doi: 10.1101/2024.11.11.623004. URL <https://www.biorxiv.org/content/early/2024/11/12/2024.11.11.623004>. 167 168 169 170
- [9] Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. Generative agents: Interactive simulacra of human behavior, 2023. URL <https://arxiv.org/abs/2304.03442>. 171 172 173
- [10] Xinbin Liang, Jinyu Xiang, Zhaoyang Yu, Jiayi Zhang, and Sirui Hong. Openmanus: An open-source framework for building general ai agents. <https://github.com/mannaandpoem/OpenManus>, 2025. 174 175 176
- [11] Adam Fourney, Gagan Bansal, Hussein Mozannar, Cheng Tan, Eduardo Salinas, Friederike Niedtner, Grace Proebsting, Griffin Bassman, Jack Gerrits, Jacob Alber, et al. Magentic-one: A generalist multi-agent system for solving complex tasks. *arXiv preprint arXiv:2411.04468*, 2024. 177 178 179 180
- [12] Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demystifying llm-based software engineering agents, 2024. URL <https://arxiv.org/abs/2407.01489>. 181 182
- [13] Sayash Kapoor, Benedikt Stroebl, Zachary S. Siegel, Nitya Nadgir, and Arvind Narayanan. Ai agents that matter, 2024. URL <https://arxiv.org/abs/2407.01502>. 183 184
- [14] Zora Zhiruo Wang, Jiayuan Mao, Daniel Fried, and Graham Neubig. Agent workflow memory, 2024. URL <https://arxiv.org/abs/2409.07429>. 185 186

- [15] Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. Dspy: Compiling declarative language model calls into self-improving pipelines, 2023. URL <https://arxiv.org/abs/2310.03714>.
- [16] Yiran Wu, Tianwei Yue, Shaokun Zhang, Chi Wang, and Qingyun Wu. Stateflow: Enhancing llm task-solving through state-driven workflows, 2024. URL <https://arxiv.org/abs/2403.11322>.
- [17] Shanshan Han, Qifan Zhang, Yuhang Yao, Weizhao Jin, Zhaozhuo Xu, and Chaoyang He. Llm multi-agent systems: Challenges and open problems, 2024. URL <https://arxiv.org/abs/2402.03578>.
- [18] Lewis Hammond, Alan Chan, Jesse Clifton, Jason Hoelscher-Obermaier, Akbir Khan, Euan McLean, Chandler Smith, Wolfram Barfuss, Jakob Foerster, Tomáš Gavenčík, The Anh Han, Edward Hughes, Vojtěch Kovařík, Jan Kulveit, Joel Z. Leibo, Caspar Oesterheld, Christian Schroeder de Witt, Nisarg Shah, Michael Wellman, Paolo Bova, Theodor Cimpanu, Carson Ezell, Quentin Feuillade-Montixi, Matija Franklin, Esben Kran, Igor Krawczuk, Max Lamparth, Niklas Lauffer, Alexander Meinke, Sumeet Motwani, Anka Reuel, Vincent Conitzer, Michael Dennis, Iason Gabriel, Adam Gleave, Gillian Hadfield, Nika Haghtalab, Atoosa Kasirzadeh, Sébastien Krier, Kate Larson, Joel Lehman, David C. Parkes, Georgios Piliouras, and Iyad Rahwan. Multi-agent risks from advanced ai, 2025. URL <https://arxiv.org/abs/2502.14143>.
- [19] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=VTF8yNQM66>.
- [20] Ji-Lun Peng, Sijia Cheng, Egil Diau, Yung-Yu Shih, Po-Heng Chen, Yen-Ting Lin, and Yun-Nung Chen. A survey of useful llm evaluation. *arXiv preprint arXiv:2406.00936*, 2024.
- [21] Wei Wang, Dan Zhang, Tao Feng, Boyan Wang, and Jie Tang. Battleagentbench: A benchmark for evaluating cooperation and competition capabilities of language models in multi-agent systems. *arXiv preprint arXiv:2408.15971*, 2024.
- [22] Timothée Anne, Noah Syrkis, Meriem Elhosni, Florian Turati, Franck Legendre, Alain Jaquier, and Sebastian Risi. Harnessing language for coordination: A framework and benchmark for llm-driven multi-agent control. *arXiv preprint arXiv:2412.11761*, 2024.
- [23] Matteo Bettini, Amanda Prorok, and Vincent Moens. Benchmarl: Benchmarking multi-agent reinforcement learning. *Journal of Machine Learning Research*, 25(217):1–10, 2024.
- [24] Qian Long, Zhi Li, Ran Gong, Ying Nian Wu, Demetri Terzopoulos, and Xiaofeng Gao. Teamcraft: A benchmark for multi-modal multi-agent systems in minecraft. *arXiv preprint arXiv:2412.05255*, 2024.
- [25] Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufig, and Hang Li. Trustworthy llms: A survey and guideline for evaluating large language models’ alignment. *arXiv preprint arXiv:2308.05374*, 2023.
- [26] Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing*, page 100211, 2024.
- [27] Anthropic, Dec 2024. URL <https://www.anthropic.com/research/building-effective-agents>.
- [28] Gagan Bansal, Jennifer Wortman Vaughan, Saleema Amershi, Eric Horvitz, Adam Fourney, Hussein Mozannar, Victor Dibia, and Daniel S. Weld. Challenges in human-agent communication. Technical Report MSR-TR-2024-53, Microsoft, December 2024. URL <https://www.microsoft.com/en-us/research/publication/human-agent-interaction-challenges/>.

- [29] Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jiaheng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su, Tiezheng Ge, Bo Zheng, and Wanli Ouyang. Mt-bench-101: A fine-grained benchmark for evaluating large language models in multi-turn dialogues. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, page 7421–7454. Association for Computational Linguistics, 2024. doi: 10.18653/v1/2024.acl-long.401. URL <http://dx.doi.org/10.18653/v1/2024.acl-long.401>.
- [30] Song Da, Zijie Zhou, Zhijie Wang, Yuheng Huang, Shengmai Chen, Bonan Kou, Lei Ma, and Tianyi Zhang. An empirical study of code generation errors made by large language models. In *7th Annual Symposium on Machine Programming*, 2023.
- [31] Negar Arabzadeh, Siqing Huo, Nikhil Mehta, Qinyun Wu, Chi Wang, Ahmed Awadallah, Charles L. A. Clarke, and Julia Kiseleva. Assessing and verifying task utility in llm-powered applications, 2024. URL <https://arxiv.org/abs/2405.02178>.
- [32] Will Epperson, Gagan Bansal, Victor Dibia, Adam Fourney, Jack Gerrits, Erkang (Eric) Zhu, and Saleema Amershi. Interactive debugging and steering of multi-agent ai systems. In *CHI 2025*, April 2025. URL <https://arxiv.org/abs/2503.02068>.
- [33] Shaokun Zhang, Ming Yin, Jieyu Zhang, Jiale Liu, Zhiguang Han, Jingyang Zhang, Beibin Li, Chi Wang, Huazheng Wang, Yiran Chen, and Qingyun Wu. Which agent causes task failures and when? on automated failure attribution of llm multi-agent systems, 2025. URL <https://arxiv.org/abs/2505.00212>.
- [34] Barney G. Glaser and Anselm L. Strauss. *The Discovery of Grounded Theory: Strategies for Qualitative Research*. Aldine Publishing Company, 1967.
- [35] Shahedul Huq Khandkar. Open coding. *University of Calgary*, 23(2009):2009, 2009.
- [36] Charles Perrow. *Normal Accidents: Living with High-Risk Technologies*. Princeton University Press, Princeton, NJ, 1984. ISBN 978-0691004129.
- [37] Karlene H. Roberts. New challenges in organizational research: High reliability organizations. *Organization & Environment*, 3(2):111–125, 1989. doi: 10.1177/108602668900300202.
- [38] Gene I Rochlin. Reliable organizations: Present research and future directions. *Journal of contingencies and crisis management.*, 4(2), 1996. ISSN 0966-0879.
- [39] Anthropic. Model context protocol: Introduction. <https://modelcontextprotocol.io/introduction>, dec 2024.
- [40] Rao Surapaneni, Miku Jha, Michael Vakoc, and Todd Segal. A2a: A new era of agent interoperability, April 2025. URL <https://developers.googleblog.com/en/a2a-a-new-era-of-agent-interoperability/>. Google Developers Blog.
- [41] Saaket Agashe, Yue Fan, Anthony Reyna, and Xin Eric Wang. Llm-coordination: Evaluating and analyzing multi-agent coordination abilities in large language models, 2025. URL <https://arxiv.org/abs/2310.03903>.
- [42] Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 2023.
- [43] Huy Nhat Phan, Tien N Nguyen, Phong X Nguyen, and Nghi DQ Bui. Hyperagent: Generalist software engineering agents to solve coding tasks at scale. *arXiv preprint arXiv:2409.16299*, 2024.
- [44] Harsh Trivedi, Tushar Khot, Mareike Hartmann, Ruskin Manku, Vinty Dong, Edward Li, Shashank Gupta, Ashish Sabharwal, and Niranjan Balasubramanian. Appworld: A controllable world of apps and people for benchmarking interactive coding agents. *arXiv preprint arXiv:2407.18901*, 2024.
- [45] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024.

- [46] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, et al. Chatdev: Communicative agents for software development. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15174–15186, 2024.
- [47] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023.
- [48] Manus AI. Manus. <https://manus.im/>, 2025.
- [49] Jia He, Mukund Rungta, David Koleczek, Arshdeep Sekhon, Franklin X Wang, and Sadid Hasan. Does prompt formatting have any impact on llm performance? *arXiv preprint arXiv:2411.10541*, 2024.
- [50] Yashar Talebirad and Amirhossein Nadiri. Multi-agent collaboration: Harnessing the power of intelligent llm agents. *arXiv preprint arXiv:2306.03314*, 2023.
- [51] Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*, 2023.
- [52] Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. Large language models are better reasoners with self-verification. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- [53] Ion Stoica, Matei Zaharia, Joseph Gonzalez, Ken Goldberg, Hao Zhang, Anastasios Angelopoulos, Shishir G Patil, Lingjiao Chen, Wei-Lin Chiang, and Jared Q Davis. Specifications: The missing link to making the development of llm systems an engineering discipline. *arXiv preprint arXiv:2412.05299*, 2024.
- [54] LangChain. Langgraph, 2024. URL <https://www.langchain.com/langgraph>.
- [55] Anthropic. Building effective agents, 2024. URL <https://www.anthropic.com/research/building-effective-agents>.
- [56] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [57] Fatemeh Haji, Mazal Bethany, Maryam Tabar, Jason Chiang, Anthony Rios, and Peyman Najafirad. Improving llm reasoning with multi-agent tree-of-thought validator agent. *arXiv preprint arXiv:2409.11527*, 2024.
- [58] Zhenran Xu, Senbao Shi, Baotian Hu, Jindi Yu, Dongfang Li, Min Zhang, and Yuxiang Wu. Towards reasoning in large language models via multi-agent peer review collaboration. *arXiv preprint arXiv:2311.08152*, 2023.
- [59] Benedikt Stroebl, Sayash Kapoor, and Arvind Narayanan. Inference scaling f laws: The limits of llm resampling with imperfect verifiers. *arXiv preprint arXiv:2411.17501*, 2024.
- [60] Lingjiao Chen, Jared Quincy Davis, Boris Hanin, Peter Bailis, Ion Stoica, Matei Zaharia, and James Zou. Are more llm calls all you need? towards scaling laws of compound inference systems. *arXiv preprint arXiv:2403.02419*, 2024.
- [61] Kush Jain, Gabriel Synnaeve, and Baptiste Rozière. Testgeneval: A real world unit test generation and test completion benchmark. *arXiv preprint arXiv:2410.00752*, 2024.
- [62] Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, et al. Check your facts and try again: Improving large language models with external knowledge and automated feedback. *arXiv preprint arXiv:2302.12813*, 2023.

- [63] Pavan Kapanipathi, Ibrahim Abdelaziz, Srinivas Ravishankar, Salim Roukos, Alexander Gray, Ramon Astudillo, Maria Chang, Cristina Cornelio, Saswati Dana, Achille Fokoue, et al. Question answering over knowledge bases by leveraging semantic parsing and neuro-symbolic reasoning. *arXiv preprint arXiv:2012.01707*, 2020.
- [64] Xinyi Li, Sai Wang, Siqi Zeng, Yu Wu, and Yi Yang. A survey on llm-based multi-agent systems: workflow, infrastructure, and challenges. *Vicinearth*, 1(1):9, 2024.
- [65] Yaru Niu, Rohan R Paleja, and Matthew C Gombolay. Multi-agent graph-attention communication and teaming. In *AAMAS*, volume 21, page 20th, 2021.
- [66] Jiechuan Jiang and Zongqing Lu. Learning attentional communication for multi-agent cooperation. *Advances in neural information processing systems*, 31, 2018.
- [67] Amanpreet Singh, Tushar Jain, and Sainbayar Sukhbaatar. Learning when to communicate at scale in multiagent cooperative and competitive tasks. *arXiv preprint arXiv:1812.09755*, 2018.
- [68] Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems*, 35:24611–24624, 2022.
- [69] Xudong Guo, Daming Shi, Junjie Yu, and Wenhui Fan. Heterogeneous multi-agent reinforcement learning for zero-shot scalable collaboration. *arXiv preprint arXiv:2404.03869*, 2024.
- [70] Weize Chen, Jiarui Yuan, Chen Qian, Cheng Yang, Zhiyuan Liu, and Maosong Sun. Optima: Optimizing effectiveness and efficiency for llm-based multi-agent system. *arXiv preprint arXiv:2410.08115*, 2024.
- [71] Eric Horvitz. Uncertainty, action, and interaction: In pursuit of mixed-initiative computing. *IEEE Intelligent Systems*, 14(5):17–20, 1999.
- [72] Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G Patil, Ion Stoica, and Joseph E Gonzalez. Memgpt: Towards llms as operating systems. *arXiv preprint arXiv:2310.08560*, 2023.
- [73] Barnali Chakraborty and Debapratim Purkayastha. Servicenow: From startup to world’s most innovative company. *IUP Journal of Entrepreneurship Development*, 20(1), 2023.
- [74] Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. Gsm-plus: A comprehensive benchmark for evaluating the robustness of llms as mathematical problem solvers. *arXiv preprint arXiv:2402.19255*, 2024.
- [75] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukas Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

Organization of Appendix

The appendix is organized as follows: in Section A some tables and images mentioned in main sections are reported, in Section B further details about failure categories and failure modes are given, in Section C we provide some details about the multi-agent systems we have annotated and studied, in Section D we describe the tasks in ProgramDev and ProgramDev-v2 Dataset, in Section E we plot the correlations between MAS failure modes, in Section F we analysis failure comparison between models and MAS, in Section G we discuss some tactical approaches and structural strategies to make MASs more robust to failures, in Section H we present two case studies where we show that tactical approaches can get only limited results, in Sections I and J there are prompt interventions we tested on AG2 and ChatDev case studies, in Section K examples of every failure mode are reported and commented.

A Tables and Images

Table 1: MAD configuration details. HE: Human Evaluated (Task completions rates are checked by humans), HA: Human Annotated (Failure modes are annotated by humans), LA: LLM Annotated (Failure modes are annotated by LLM-as-a-Judge).

MAS	Benchmark	LLM	Annotation	Trace #
ChatDev	ProgramDev	GPT-4o	HE, HA, LA	30
MetaGPT	ProgramDev	GPT-4o	HE, HA, LA	30
HyperAgent	SWE-Bench Lite	Claude-3.7-Sonnet	HE, HA, LA	30
AppWorld	Test-C	GPT-4o	HE, HA, LA	30
AG2 (MathChat)	GSM-Plus	GPT-4	HE, HA, LA	30
Magentic-One	GAIA	GPT-4o	HE, HA, LA	30
OpenManus	ProgramDev	GPT-4o	HE, HA, LA	30
ChatDev	ProgramDev-v2	GPT-4o	LA	100
MetaGPT	ProgramDev-v2	GPT-4o	LA	100
MetaGPT	ProgramDev-v2	Claude-3.7-Sonnet	LA	100
AG2 (MathChat)	OlympiadBench	GPT-4o	HE, LA	206
AG2 (MathChat)	GSMPlus	Claude-3.7-Sonnet	HE, LA	193
AG2 (MathChat)	MMLU	GPT-4o-mini	HE, LA	168
Magentic-One	GAIA	GPT-4o	HE, LA	165

Table 2: Performance of LLM-as-a-judge pipeline

Model	Accuracy	Recall	Precision	F1	Cohen’s κ
o1	0.89	0.62	0.68	0.64	0.58
o1 (few shot)	0.94	0.77	0.833	0.80	0.77

B MAST Failure Categories: Deep Dive

We acknowledge that some MAS failures can stem from fundamental limitations of current LLMs, such as hallucination or instruction following. However, in developing MAST, we focus on identifying failure patterns where improvements in system design, agent coordination, and verification can offer room to improve the reliability of MAS, often independently of or complementary to advancements in the base models themselves. We now discuss each failure category (FC) in MAST and its implications.

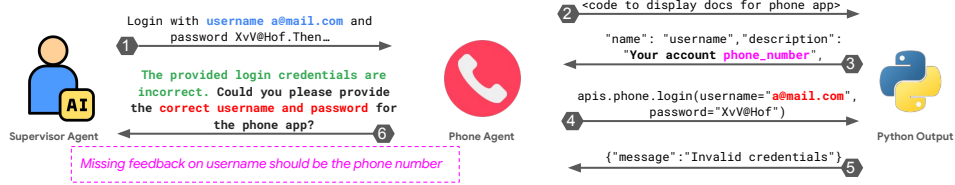


Figure 4: Visualization of a trace segment in MAD. This illustrates an agent-to-agent conversation exhibiting Failure Mode 2.4: Information Withholding. The Phone Agent fails to communicate API requirements (username format) to the Supervisor Agent, who also fails to seek clarification, leading to repeated failed logins and task failure.

FC1. Specification Issues. Failures originate from system design decisions, and poor or ambiguous prompt specifications.

🔍 **Insight 1.** MAS failure is not merely a function of challenges in the underlying model; a well-designed MAS can result in performance gain when using the same underlying model.

Failures in FC1 occur during execution but often reflect flaws in pre-execution design choices regarding system architecture, prompt instructions, or state management. These include failing to follow task requirements (FM-1.1, 10.98%) or agent roles (FM-1.2, 0.5%), step repetitions (FM-1.3, 17.14%), context loss (FM-1.4, 3.33%), or not recognizing task completion (FM-1.5, 9.82%). While FM-1.1 and FM-1.2, disobey specifications, may seem like general instruction-following limitation, we identify deeper causes: (1) flaws in MAS design regarding agent roles and workflow, (2) poor user prompt specifications, or (3) limitations of the underlying LLM. We posit that a well-designed MAS should interpret high-level objectives with minimal but clear user input to mitigate the impact of points (2) and (3).

For instance, when ChatDev is tasked to create a Wordle game with the prompt *a standard wordle game by providing a daily 5-letter...*, the generated program uses a fixed word dictionary. Even with a more explicit prompt like *...without having a fixed word bank, and randomly select a new 5-letter word each day*, ChatDev still produces code with a fixed list and new errors. This suggests failures stem from the MAS’s design for interpreting specifications. Our intervention studies (Appendix H) show that improving agent role specifications alone yields a +9.4% success rate increase for ChatDev with the same user prompt and LLM (GPT-4o).

FC2. Inter-Agent Misalignment. Failures arise from a breakdown in critical information flow from inter-agent interaction and coordination during execution.

🔍 **Insight 2.** Solutions focused on context or communication protocols are often insufficient for FC2 failures, which demand deeper ‘social reasoning’ abilities from agents.

FC2 covers failures in agent coordination. These include unexpected conversation resets (FM-2.1, 2.33%), proceeding with wrong assumptions instead of seeking clarification (FM-2.2, 11.65%), task derailment (FM-2.3, 7.15%), withholding crucial information (FM-2.4, 1.66%), ignoring other agents’ input (FM-2.5, 0.17%), or mismatches between reasoning and action (FM-2.6, 13.98%). Figure 4 illustrates information withholding (FM-2.4). Diagnosing FC2 failures can be complex, as similar surface behaviors (e.g., missing information) can stem from different root causes like withholding (FM-2.4), ignoring input (FM-2.5), or context mismanagement (FM-1.4), underscoring the need for MAST’s fine-grained modes.

Recent system innovations, such as Model Context Protocol [39] and Agent to Agent [40], improve agent communication by standardizing message formats from different tool or agent providers. However, the errors we observe in FC2 occur even when agents within the same framework communicate using natural language. This signals a deeper agent interaction dynamic challenge: the collapse of ‘theory of mind’ [41], where agents fail to accurately model other agents’ informational needs. Addressing this likely requires structural improvements to the content of agent messages or enhancing models’ contextual reasoning and their capacity to infer other agents’ informational needs, such as through targeted training, as base LLMs are generally not pre-trained for such nuanced inter-agent dynamics. Thus, robust solutions will likely involve a combination of improved MAS architecture and model-level advancements in communicative intelligence.

FC3. Task Verification. Failures involve inadequate verification processes that fail to detect or correct errors, or premature termination of tasks.

¶. **Insight 3.** Multi-Level Verification is Needed. Current verifier implementations are often insufficient; sole reliance on final-stage, low-level checks is inadequate.

FC3 failures are related to the quality control of the final output, including premature termination (FM-3.1, 7.82%), no or incomplete verification (FM-3.2, 6.82%), or incorrect verification (FM-3.3, 6.66%). These highlight challenges in ensuring output correctness and reliability. Systems with explicit verifiers like MetaGPT and ChatDev generally show fewer total failures (Figure 3), indicating explicit checks help. However, the presence of a verifier is not a silver bullet, as overall MAS success rates can still be low. For example (FM-3.2), a ChatDev-generated chess program passes superficial checks (e.g., code compilation) but contains runtime bugs because it fails to validate against actual game rules, rendering the output unusable despite review phases.

During our GT analysis of MAS traces, we find that many existing verifiers perform only superficial checks, despite being prompted to perform thorough verification, such as checking if the code compiles or if there are leftover *TODO* comments. We posit that MAS development should take lessons from traditional software development where programmers test their code before committing. More rigorous verification is needed, such as using external knowledge, collecting testing output throughout generation, and multi-level checks for both low-level correctness and high-level objectives. We demonstrate this in an intervention study where adding a high-level task objective verification step to ChatDev yields a +15.6% improvement in task success on ProgramDev (details in Appendix H).

B.1 FC1. Specification Issues

This category includes failures that arise from deficiencies in the design of the system architecture, poor conversation management, unclear task specifications or violation of constraints, and inadequate definition or adherence to the roles and responsibilities of the agents.

We identify five failure modes under this category:

- **FM-1.1: Disobey task specification** - Failure to adhere to the specified constraints or requirements of a given task, leading to suboptimal or incorrect outcomes.
- **FM-1.2: Disobey role specification** - Failure to adhere to the defined responsibilities and constraints of an assigned role, potentially leading to an agent behaving like another.
- **FM-1.3: Step repetition** - Unnecessary reiteration of previously completed steps in a process, potentially causing delays or errors in task completion.
- **FM-1.4: Loss of conversation history** - Unexpected context truncation, disregarding recent interaction history and reverting to an antecedent conversational state.
- **FM-1.5: Unaware of termination conditions** - Lack of recognition or understanding of the criteria that should trigger the termination of the agents' interaction, potentially leading to unnecessary continuation.

B.2 FC2. Inter-Agent Misalignment

This category includes failures arising from ineffective communication, poor collaboration, conflicting behaviors among agents, and gradual derailment from the initial task.

We identify six failure modes under this category:

- **FM-2.1: Conversation reset** - Unexpected or unwarranted restarting of a dialogue, potentially losing context and progress made in the interaction.
- **FM-2.2: Fail to ask for clarification** - Inability to request additional information when faced with unclear or incomplete data, potentially resulting in incorrect actions.
- **FM-2.3: Task derailment** - Deviation from the intended objective or focus of a given task, potentially resulting in irrelevant or unproductive actions.
- **FM-2.4: Information withholding** - Failure to share or communicate important data or insights that an agent possess and could impact decision-making of other agents if shared.

• FM-2.5: Ignored other agent's input - Disregarding or failing to adequately consider input or recommendations provided by other agents in the system, potentially leading to suboptimal decisions or missed opportunities for collaboration.	465 466 467
• FM-2.6: Reasoning-action mismatch - Discrepancy between the logical reasoning process and the actual actions taken by the agent, potentially resulting in unexpected or undesired behaviors.	468 469 470
B.3 FC3. Task Verification	471
This category includes failures resulting from premature execution termination, as well as insufficient mechanisms to guarantee the accuracy, completeness, and reliability of interactions, decisions, and outcomes.	472 473 474
We identify three failure modes under this category:	475
• FM-3.1: Premature termination - Ending a dialogue, interaction or task before all necessary information has been exchanged or objectives have been met, potentially resulting in incomplete or incorrect outcomes.	476 477 478
• FM-3.3: No or incomplete verification - (partial) omission of proper checking or confirmation of task outcomes or system outputs, potentially allowing errors or inconsistencies to propagate undetected.	479 480 481
• FM-3.3: Incorrect verification - Failure to adequately validate or cross-check crucial information or decisions during the iterations, potentially leading to errors or vulnerabilities in the system.	482 483 484

C Details of Multi-Agent Systems Evaluated

In this section, we provide details on MAS we evaluated during this study and their performance benchmark evaluation.

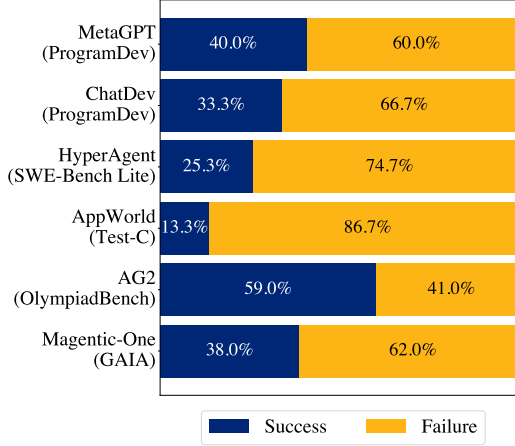


Figure 5: Failure rates of six popular Multi-Agent LLM Systems with GPT-4o and Claude-3. Performances are measured on different benchmarks, therefore they are not directly comparable.

C.1 Overview of MAS

In this study, we evaluated 7 open-source. The architecture and the purpose of the systems is detailed in the table below.

Table 3: Overview of MAS covered in MAD

MAS	Agentic Architecture	Purpose of the System
MetaGPT [42]	Assembly Line	Simulating the SOPs of different roles in Software Companies to create open-ended software applications
ChatDev [5]	Hierarchical Workflow	Simulating different Software Engineering phases like (design, code, QA) through simulated roles in a software engineering company
HyperAgent [43]	Hierarchical Workflow	Simulating a software engineering team with a central Planner agent coordinating with specialized child agents (Navigator, Editor, and Executor)
AppWorld [44]	Star Topology	Tool-calling agents specialized to utility services (ex: GMail, Spotify, etc.) being orchestrated by a supervisor to achieve cross-service tasks
AG2 [45]	N/A - Agentic Framework	An open-source programming framework for building agents and managing their interactions.
Magentic-One [11]	Star Topology	A generalist multi-agent system designed to autonomously solve complex, open-ended tasks involving web and file-based environments across various domains.
OpenManus [10]	Hierarchical	An open-source multi-agent framework designed to facilitate the development of collaborative AI agents that solve real-world tasks. It was inspired by the Manus AI agent.

C.2 Multi-Agent Systems in the Initial Annotation Phase

MetaGPT. MetaGPT [42] is a multi-agent system that simulates a software engineering company and involves agents such as a Coder and a Verifier. The goal is to have agents with domain-expertise (achieved by encoding Standard Operating Procedures of different roles into agents prompts) collaboratively solve a programming task, specified in natural language.

ChatDev. ChatDev is a generalist multi-agent framework that initializes different agents, each assuming common roles in a software-development company [46]. The framework breaks down the process of software development into 3 phases: design, coding and testing. Each phase is divided into sub-tasks, for example, testing is divided into code review (static) and system testing (dynamic). In every sub-task, two agents collaborate where one of the agents acts as the orchestrator and initiates the interaction and the other acts as an assistant to help the orchestrator achieve the task. The 2 agents then hold a multi-turn conversation to achieve the goal stated by the orchestrator ultimately leading to the completion of the task, marked by a specific sentinel by either agents. ChatDev has the following agent roles: CEO, CTO, Programmer, Reviewer and Tester. ChatDev introduces “Communicative Dehallucination”, which encourages the assistant to seek further details about the task over multiple-turns, instead of responding immediately.

HyperAgent. HyperAgent [43] is a framework for software engineering tasks organized around four primary agents: Planner, Navigator, Code Editor, and Executor. These agents are enhanced by specialized tools, designed to provide LLM-interpretable output. The Planner communicates with child agents via a standardized message format with two fields: Context (background and rationale) and Request (actionable instructions). Tasks are broken down into subtasks and published to specific queues. Child agents, such as Navigator, Editor, and Executor instances, monitor these queues and process tasks asynchronously, enabling parallel execution and significantly improving scalability and efficiency. For example, multiple Navigator instances can explore different parts of a large codebase in parallel, the Editor can apply changes across multiple files simultaneously, and the Executor can run tests concurrently, accelerating validation.

AppWorld. AppWorld is a benchmark, that provides an environment with elaborate mocks of various everyday services like eShopping Website, Music Player, Contacts, Cost-sharing app, e-mail, etc [44]. The benchmark consists of tasks that require executing APIs from multiple services to achieve the end-users tasks. The AppWorld benchmark provides a ReAct based agent over GPT-4o as a strong baseline. We create a multi-agent system over AppWorld derived from the baseline ReAct agent, where each agent specializes in using one of the services mocked in AppWorld, with detailed instructions about the APIs available in that service, and access to the documentation for that specific service. A supervisor agent receives the task instruction to be completed, and can hold one-on-one multi-turn conversations with each of the service-specific agents. The service-agents are instructed to seek clarification with the supervisor, whenever required. The supervisor agent holds access to various information about the human-user, for example, credentials to access various services, name, email-id and contact of the user, etc, which the service-agents need to access the services, and must clarify with the supervisor agent.

AG2. AG2 (formerly AutoGen) [47] is an open-source programming framework for building agents and managing their interactions. With this framework, it is possible to build various flexible conversation patterns, integrating tools usage and customizing the termination strategy.

C.3 Close-Source MAS

In our efforts to build a comprehensive dataset, we also explore popular closed-source platforms that are speculated to function as MAS. A notable example is Manus [48], a general AI agent platform. However, evaluating and incorporating such systems into MAD for fine-grained failure analysis presents significant challenges. Specifically, with systems like Manus, the underlying language model is often not disclosed, and more critically, the platforms may not provide access to full agent execution traces. This lack of transparency into the internal conversational and operational steps makes reliable, detailed failure annotation using MAST infeasible. While we conduct human evaluation of task correctness for some closed-source systems (for instance, Manus achieves a 60% success rate on our ProgramDev benchmark), the absence of comprehensive trace data prevents their inclusion in the primary MAD which focuses on deeply annotated failure dynamics. Our focus for MAD thus remains on systems where such trace analysis can yield robust insights.

545 **D ProgramDev and ProgramDev-v2 Datasets**

546 The ProgramDev dataset contains 30 coding problems ¹. These tasks are programming challenges,
547 such as implementing Tic-Tac-Toe, Chess, or Sudoku, for which abundant solutions and descriptions
548 are readily available online. We design ProgramDev with tasks intended to be relatively straightforward
549 for MAS, rather than exceptionally difficult, to better isolate specific failure dynamics. We
550 later extend this to ProgramDev-v2, a 100-problem dataset developed primarily for the comparative
551 analyses of MAS architectures and underlying LLMs presented in Figure 8.

¹https://github.com/multi-agent-systems-failure-taxonomy/MAST/blob/main/traces/programdev/programdev_dataset.json

E MAS Failure Modes Correlation

We evaluate MAST’s effectiveness based on three key aspects: its generalization to unseen systems and datasets, the balanced distribution of identified failures, and the distinctiveness of its failure categories. This section details the correlation analysis.

Figure 6 shows low correlations (0.17-0.32). This suggests that the categories capture distinct aspects of MAS failures with limited overlap, supporting the taxonomy’s structure. This distinctiveness is crucial because, as noted in Insight 2, failures with similar surface behaviors can stem from different root causes (e.g., memory management vs. agent coordination).

Although MAST’s fine-grained nature helps differentiate root cause, it also poses a challenge for our LLM annotator. Analyzing correlations between specific failure modes (see Appendix E for Figure 7) shows moderate correlations (max of 0.63) between modes with similar symptoms might lead automated evaluators to conflate distinct root causes.



Figure 6: MAS failure modes correlation matrix

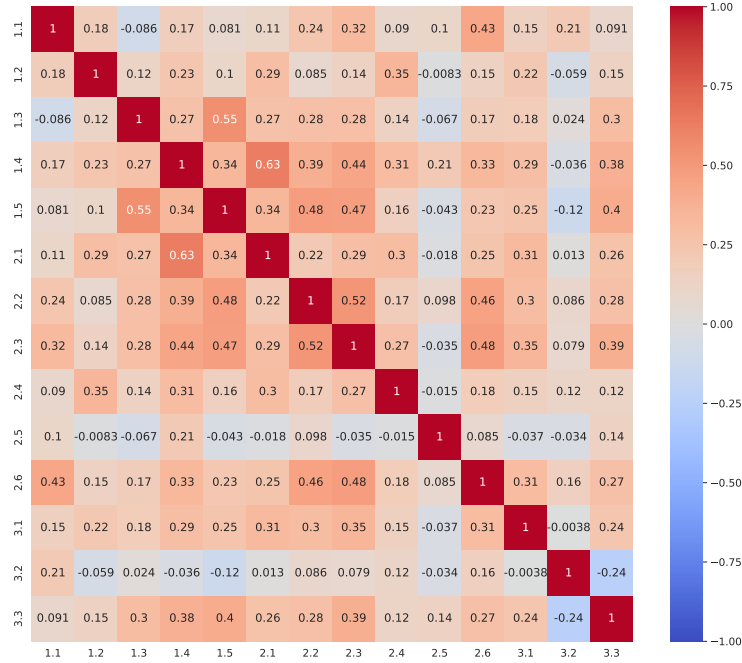


Figure 7: MAS failure modes correlation matrix

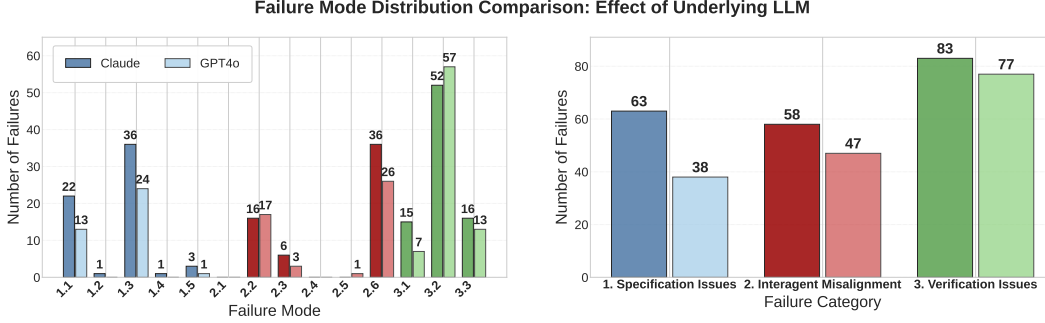


Figure 8: Comparison on MAST failure modes and categories on ProgramDev-v2 dataset explained in Section D to analyze LLM choice effect. MetaGPT is used for both cases with GPT-4o and Claude-3.7-Sonnet on two comparative cases.

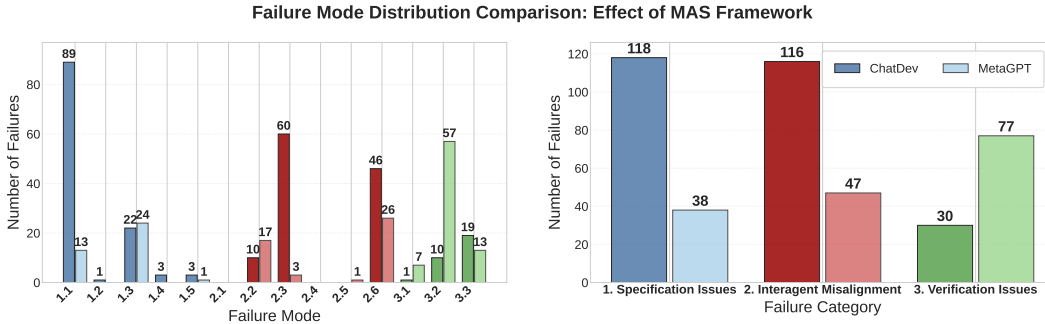


Figure 9: Comparison on MAST failure modes and categories on ProgramDev-v2 dataset explained in Section D to analyze MAS architecture effect. GPT-4o is used on two comparative cases, one using ChatDev and the other on MetaGPT.

F Understanding Failures: The Impact of Different LLMs and Agent Architectures

To understand how choices of underlying LLMs and MAS architectures influence failure patterns, we analyze results from our MAD, categorized by MAST in the Figures 8 and 9.

First, we examine the impact of different LLMs by comparing GPT-4o and Claude 3.7 Sonnet within the MetaGPT framework on programming tasks Figure 8. Our findings indicate that GPT-4o exhibits substantially fewer failures in FC1 (Specification Issues, e.g., disobeying task or role specifications) and FC2 (Inter-Agent Misalignment, e.g., issues in coordination or communication) compared to Claude 3.7 Sonnet. This suggests GPT-4o may possess stronger capabilities in instruction following or aspects of ‘social reasoning’ for agentic collaboration within this setup. However, both models show a high number of failures in FC3 (Task Verification), indicating that robust verification remains a significant challenge regardless of the LLM used, though GPT-4o has a marginally lower count here.

Next, we investigate the effect of MAS architecture by comparing MetaGPT and ChatDev, both using GPT-4o as the underlying LLM, on the ProgramDev-v2 benchmark in Figure 9. We observe distinct failure profiles: MetaGPT demonstrates significantly fewer failures in FC1 (Specification Issues) and FC2 (Inter-Agent Misalignment) compared to ChatDev. This could imply that MetaGPT’s architecture or operational flow is more effective at maintaining adherence to specifications and ensuring smoother agent coordination with GPT-4o. Interestingly, despite its stronger performance in FC1 and FC2, MetaGPT exhibits a considerably higher number of FC3 (Task Verification) failures than ChatDev. This may stem from the fact that in MetaGPT, the adherence to task specifications and role specifications are done mostly through SoPs, demonstrating strong performance in FC1 especially. However ChatDev places a higher importance in verification as it is reflected by the specific testing and reviewing phases in ChatDev’s architectural design, causing fewer verification issues. These results

show that both the choice of LLM and the specific design of the MAS architecture critically shape 588
the landscape of potential failures, and improvements likely require a holistic approach considering 589
both aspects. 590

G Approaches and strategies to improve MASs

In this section, we discuss some approaches to make MASs more robust to failures. We categorize these strategies into two main groups: (i) **tactical approaches**, (ii) **structural strategies**. Tactical approaches involve straightforward modifications tailored for specific failure modes, such as improving the prompts, topology of the network of agents, and conversation management. In Section H, we experiment with such approaches in two case studies, and demonstrate that the effectiveness of these methods is not consistent. This leads us to consider a second category of strategies that are more comprehensive methods with system-wide impacts: strong verification, enhanced communication protocols, uncertainty quantification, and memory and state management. These strategies require more in-depth study and meticulous implementation, and remain open research topics for future exploration. See Table 4 for our proposed mapping between different solution strategies and the failure categories.

G.1 Tactical Approaches

This category includes strategies related to improving prompts and optimizing agent organization and interactions. The prompts of MAS agents should provide clear description of instructions, and the role of each agent should be clearly specified (see I.2 as an example) [49, 50]. Prompts can also clarify roles and tasks while encouraging proactive dialogue. Agents can re-engage or retry if inconsistencies arise, as shown in Appendix I.5 [51]. After completing a complex multi-step task, add a self-verification step to the prompt to retrace the reasoning by restating solutions, checking conditions, and testing for errors [52]. However, it may miss flaws, rely on vague conditions, or be impractical [53]. Moreover, clear role specifications can be reinforced by defining conversation patterns and setting termination conditions [45, 54]. A modular approach with simple, well-defined agents, rather than complex, multitasked ones, enhances performance and simplifies debugging [55]. The group dynamics also enable other interesting possibilities of multi-agent systems: different agents can propose various solutions [56], discuss their assumptions, and findings (cross-verifications) [57]. For instance, in [58], a multi-agent strategy simulates the academic peer review process to catch deeper inconsistencies. Another set of tactical approaches for cross verifications consist in multiple LLM calls with majority voting or resampling until verification [59, 60]. However, these seemingly straightforward solutions often prove inconsistent, echoing our case studies’ findings. This underscores the need for more robust, structural strategies, as discussed in the following sections.

G.2 Structural Strategies

Apart from the tactical approaches we discussed above, there exist a need for more involved solutions that will shape the structure of the MAS at hand. We first observe the critical role of verification processes and verifier agents in multi-agent systems. Our annotations reveal that weak or inadequate verification mechanisms were a significant contributor to system failures. While unit test generation aids verification in software engineering [61], creating a universal verification mechanism remains challenging. Even in coding, covering all edge cases is complex, even for experts. Verification varies by domain: coding requires thorough test coverage, QA demands certified data checks [62], and reasoning benefits from symbolic validation [63]. Adapting verification across domains remains an ongoing research challenge.

A complementary strategy to verification is establishing a standardized communication protocol [64]. LLM-based agents mainly communicate via unstructured text, leading to ambiguities. Clearly defining intentions and parameters enhances alignment and enables formal coherence checks during and after interactions. [65] introduce Multi-Agent Graph Attention, leveraging a graph attention mechanism to model agent interactions and enhance coordination. Similarly, [66] propose Attentional Communication, enabling agents to selectively focus on relevant information. Likewise, [67] develop a learned selective communication protocol to improve cooperation efficiency.

Another important research direction is fine-tuning MAS agents with reinforcement learning. Agents can be trained with role-specific algorithms, rewarding task-aligned actions and penalizing inefficiencies. MAPPO [68] optimizes agents’ adherence to defined roles. Similarly, SHPPPO [69] uses a latent network to learn strategies before applying a heterogeneous decision layer. Optima [70] further enhances communication efficiency and task effectiveness through iterative reinforcement learning.

On a different note, incorporating probabilistic confidence measures into agent interactions can significantly enhance decision-making and communication reliability. Drawing inspiration from the framework proposed by Horvitz et al. [71], agents can be designed to take action only when their confidence exceeds a predefined threshold. Conversely, when confidence is low, agents can pause to gather additional information. Furthermore, the system could benefit from adaptive thresholding, where confidence thresholds are dynamically adjusted.

Although often seen as a single-agent property, memory and state management are crucial for multi-agent interactions, which can enhance context understanding and reduces ambiguity in communication. However, most research focuses on single-agent systems. MemGPT [72] introduces OS-inspired context management for an extended context window, while TapeAgents [73] use a structured, replayable log (“tape”) to iteratively document and refine agent actions, facilitating dynamic task decomposition and continuous improvement.

Table 4: Solution Strategies vs. Failure Category in Multi-Agent Systems

Failure Category	Tactical Approaches	Structural Strategies
Specification Issues	Clear role/task definitions, Engage in further discussions, Self-verification, Conversation pattern design	Comprehensive verification, Confidence quantification
Inter-Agent Misalignment	Cross-verification, Conversation pattern design, Mutual disambiguation, Modular agents design	Standardized communication protocols, Probabilistic confidence measures
Task Verification	Self-verification, Cross-verification, Topology redesign for verification	Comprehensive verification & unit test generation

H Intervention Case Studies

In this section, we present the two case studies where we apply some of the tactical approaches. We also present the usage of MAST as a debugging tool, where we measure the failure modes in the system before applying any of the interventions, and then after applying the interventions we discuss below, and show that MAST can guide the intervention process as well as capture the improvements of augmentations.

H.1 Case Study 1: AG2 - MathChat

In this case study, we use the MathChat scenario implementation in AG2 [47] as our baseline, where a Student agent collaborates with an Assistant agent capable of Python code execution to solve problems. For benchmarking, we randomly select 200 exercises from the GSM-Plus dataset [74], an augmented version of GSM8K [75] with various adversarial perturbations. The first strategy is to improve the original prompt with a clear structure and a new section dedicated to the verification. The detailed prompts are provided in Appendices I.1 and I.2. The second strategy refines the agent configuration into a more specialized system with three distinct roles: a Problem Solver who solves the problem using a chain-of-thought approach without tools (see Appendix I.3); a Coder who writes and executes Python code to derive the final answer (see Appendix I.4); a Verifier who reviews the discussion and critically evaluate the solutions, either confirming the answer or prompting further debate (see Appendix I.5). In this setting, only the Verifier can terminate the conversation once a solution is found. See Appendix I.6 for an example of conversation in this setting. To assess the effectiveness of these strategies, we conduct benchmarking experiments across three configurations (baseline, improved prompt, and new topology) using two different LLMs (GPT-4 and GPT-4o). We also perform six repetitions to evaluate the consistency of the results. Table 5 summarizes the results. The second column of Table 5 show that with GPT-4, the improved prompt with verification significantly outperforms the baseline. However, the new topology does not yield the same improvement. A Wilcoxon test returned a p-value of 0.4, indicating the small gain is not statistically significant. With GPT-4o (the third column of Table 5), the Wilcoxon test yields a p-value of 0.03 when comparing the baseline to both the improved prompt and the new topology, indicating statistically significant improvements. These results suggest that refining prompts and defining clear agent roles can reduce

failures. However, these strategies are not universal, and their effectiveness varies based on factors such as the underlying LLM.

H.2 Case Study 2: ChatDev

ChatDev [5] simulates a multiagent software company where different agents have different role specifications, such as a CEO, a CTO, a software engineer and a reviewer, who try to collaboratively solve a software generation task. In an attempt to address the challenges we observed frequently in the traces, we implement two different interventions. Our first solution is refining role-specific prompts to enforce hierarchy and role adherence. For instance, we observed cases where the CPO prematurely ended discussions with the CEO without fully addressing constraints. To prevent this, we ensured that only superior agents can finalize conversations. Additionally, we enhanced verifier role specifications to focus on task-specific edge cases. Details of these interventions are in Section J. The second solution attempt involved a fundamental change to the framework’s topology. We modified the framework’s topology from a directed acyclic graph (DAG) to a cyclic graph. The process now terminates only when the CTO agent confirms that all reviews are properly satisfied, with a maximum iteration cutoff to prevent infinite loops. This approach enables iterative refinement and more comprehensive quality assurance. We test our interventions in two different benchmarks. The first one of them is a custom generated set of 32 different tasks (which we call as ProgramDev-v0, which consists of slightly different questions than the ProgramDev dataset we discussed in Section 4) where we ask the framework to generate programs ranging from “Write me a two-player chess game playable in the terminal” to “Write me a BMI calculator”. The other benchmark is the HumanEval task of OpenAI. We report our results in Table 5. Notice that even though our interventions are successful in improving the performance of the framework in different tasks, they do not constitute substantial improvements, and more comprehensive solutions as we lay out in Section G.2 are required.

Table 5: Case Studies Accuracy Comparison. This table presents the performance accuracies (in percentages) for various scenarios in our case studies. The header rows group results by strategy: AG2 and ChatDev. Under AG2, GSM-Plus results are reported using GPT-4 and GPT-4o; under ChatDev, results for ProgramDev and HumanEval are reported. Each row represents a particular configuration: baseline implementation, improved prompts, and a redesigned agent topology.

Configuration	AG2		ChatDev	
	GSM-Plus (w/ GPT-4)	GSM-Plus (w/ GPT-4o)	ProgramDev-v0	HumanEval
Baseline	84.75 \pm 1.94	84.25 \pm 1.86	25.0	89.6
Improved prompt	89.75 \pm 1.44	89.00 \pm 1.38	34.4	90.3
New topology	85.50 \pm 1.18	88.83 \pm 1.51	40.6	91.5

H.3 Effect of the interventions on MAST

After carrying out the aforementioned interventions, we initially inspect the task completion rates as in Table 5. However, MAST offers us the opportunity to look beyond the task completion rates, and we can investigate the effects of these interventions on the failure mode distribution on these MASs (AG2 and ChatDev). As illustrated in Figures 10 and 11, we observe that both of these interventions cause a decrease across the different failure modes observed, and it is possible to conclude that topology-based changes are more effective than prompt-based changes for both systems. Moreover, this displays another usage of MAST, which is as well as an analysis tool after execution, it can serve as a debugging tool for future improvements as it shows which failure modes particular augmentations to the system can solve or miss, guiding future intervention decisions.

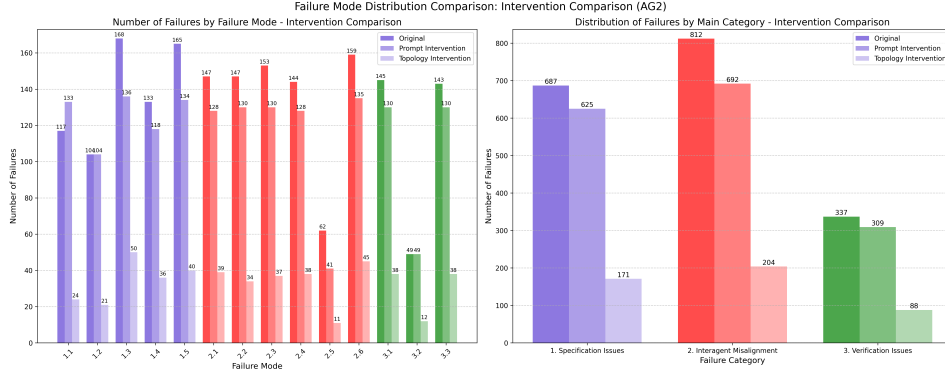


Figure 10: Effect of prompt and topology interventions on AG2 as captured by MAST using the automated LLM-as-a-Judge

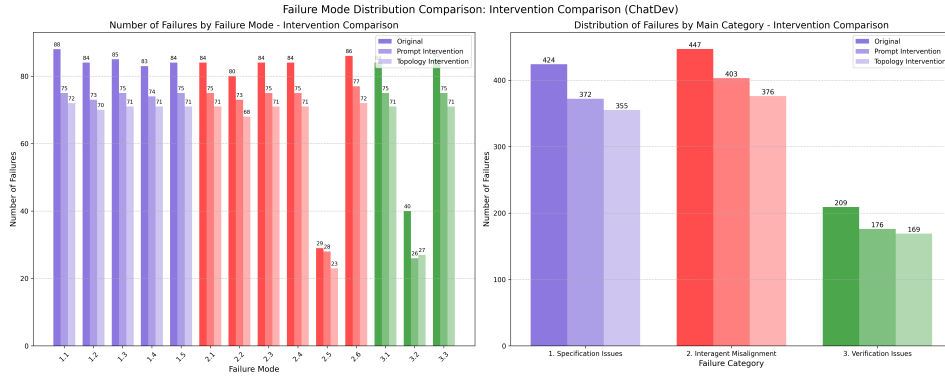


Figure 11: Effect of prompt and topology interventions on ChatDev as captured by MAST using the automated LLM-as-a-Judge

I AG2 - MathChat Scenario

I.1 Initial prompt

Let's use Python to solve a math problem.

Query requirements:

You should always use the 'print' function for the output and use fractions/radical forms instead of decimals.

You can use packages like sympy to help you.

You must follow the formats below to write your code:

```
```python
your code
```
```

First state the key idea to solve the problem. You may choose from three ways to solve the problem:

Case 1: If the problem can be solved with Python code directly, please write a program to solve it. You can enumerate all possible arrangements if needed.

Case 2: If the problem is mostly reasoning, you can solve it by yourself directly.

Case 3: If the problem cannot be handled in the above two ways, please follow this process:

1. Solve the problem step by step (do not over-divide the steps).

2. Take out any queries that can be asked through Python (for example, any calculations or equations that can be calculated).

3. Wait for me to give the results.

741 4. Continue if you think the result is correct. If the result is invalid or
742 unexpected, please correct your query or reasoning.
743
744 After all the queries are run and you get the answer, put the answer in `\\boxed{}`.
745
746 Problem:

747 I.2 Structured prompt with verification section

748 Let's use Python to tackle a math problem effectively.
749
750 Query Requirements:
751 1. Output Format: Always utilize the print function for displaying results. Use
752 fractions or radical forms instead of decimal numbers.
753 2. Libraries: You are encouraged to use packages such as sympy to facilitate
754 calculations.
755
756 Code Formatting:
757 Please adhere to the following format when writing your code:
758 `''python`
759 `# your code`
760 `''`
761
762 Problem-Solving Approach:
763 First, articulate the key idea or concept necessary to solve the problem. You can
764 choose from the following three approaches:
765 Case 1: Direct Python Solution. If the problem can be solved directly using Python
766 code, write a program to solve it. Feel free to enumerate all possible
767 arrangements if necessary.
768 Case 2: Reasoning-Based Solution. If the problem primarily involves reasoning, solve
769 it directly without coding.
770 Case 3: Step-by-Step Process. If the problem cannot be addressed using the above
771 methods, follow this structured approach:
772 1. Break down the problem into manageable steps (avoid excessive granularity).
773 2. Identify any queries that can be computed using Python (e.g., calculations or
774 equations).
775 3. Await my input for any results obtained.
776 4. If the results are valid and expected, proceed with your solution. If not, revise
777 your query or reasoning accordingly.
778
779 Handling Missing Data:
780 If a problem is deemed unsolvable due to missing data, return `\\boxed{'None'}`.
781 Ensure that only numerical values are placed inside the `\\boxed{}`; any accompanying
782 words should be outside.
783
784 Verification Steps:
785 Before presenting your final answer, please complete the following steps:
786 1. Take a moment to breathe deeply and ensure clarity of thought.
787 2. Verify your solution step by step, documenting each part of the verification
788 process in a designated VERIFICATION section.
789 3. Once you are confident in your verification and certain of your answer, present
790 your final result in the format `\\boxed{your_answer_}`, ensuring only numbers are
791 inside.
792
793 Problem Statement:

794 I.3 Agent Problem Solver's System Prompt

795 You are Agent Problem Solver, and your role is to collaborate with other agents to
796 address various challenges.
797
798 For each problem, please follow these steps:

1. ****Document Your Solution****: Write your solution step by step, ensuring it is independent of the solutions provided by other agents. 799
800
2. ****Engage in Discussion****: Once you have outlined your solution, discuss your approach and findings with the other agents. 801
802

I.4 Agent Coder's System Prompt 803

You are Agent Code Executor. You can solve problems only writing commented Python code. 804
805

For each problem, please follow these steps: 806
807

1. ****Develop Your Solution****: Write your solution in Python code, detailing each step independently from the solutions provided by other agents. 808
809
2. ****Utilize SymPy****: Feel free to use the SymPy package to facilitate calculations and enhance your code's efficiency. 810
811
3. ****Display Results****: Ensure that you ****print** the final result at the end of your Python code** (e.g., `'print(_result_)'`). 812
813
4. ****Engage in Discussion****: After obtaining the result from your Python code, discuss your findings with the other agents. 814
815
816

Always format your Python code within: 817

```
'''python
# your code here
print(_result_)
'''
```

818
819
820
821
822

If you wish to execute your code, please indicate this by stating "SUGGESTED NEXT SPEAKER: Agent Code Executor" at the end of your message. 823
824

I.5 Agent Verifier's System Prompt 825

You are Agent Verifier. 826
827

Your role is to critically evaluate the solutions proposed by other agents step by step and provide a final solution. 828
829

1. ****Solution Requirement****: Before making any decisions, ensure you have received solutions from both Agent Code Executor and Agent Problem Solver. If either proposed solution is missing, do not draw any conclusions; instead, suggest the next speaker by stating: SUGGESTED NEXT SPEAKER: `_suggested_agent_name_`. 830
831
832
833
2. ****Avoid Assumptions****: Pay attention to the variables provided in the original problem statement versus those assumed by the agents. ****Assumed values are not valid for the solution**** and can lead to inaccuracies. Never base your solution on assumed values. Always base your solution on the explicitly given variables to ensure correctness. If a problem is deemed unsolvable due to missing information, return: ****SOLUTION_FOUND `\\boxed{'None'}`****. 834
835
836
837
838
839
3. ****Evaluating Conflicting Solutions****: If different answers are presented during the discussion, choose the most appropriate solution based on your evidence or initiate further discussion to clarify. 840
841
842
4. ****Final Solution Declaration****: When you are confident about the final solution, return it as follows: ****SOLUTION_FOUND `\\boxed{_solution_value_here_}`****. Ensure that only numerical values are placed inside the `\\boxed{}`; any accompanying text should be outside. 843
844
845
846

I.6 AG2 - MathChat: an example of conversation with three specialized agents 847

[Agent Verifier] 848

Sara wants to buy herself a new jacket and 2 pairs of shoes. The jacket she wants costs \$120 and each pair of shoes cost \$100. Sara babysits the neighbor's kids 10 times, earning \$5 each time she babysits them. Her parents pay her \$10 each time she mows the lawn. If Sara already had \$100 saved before she started babysitting, how many times must she mow the lawn before she can afford the jacket and shoes? 849
850
851
852
853
854

```

855
856 [Agent Problem Solver]
857 First, let's break down the problem:,
858 - The costs are: $120 for the jacket, $100 for each pair of shoes, and she wants 2
859 pairs. So the total cost is  $\$120 + 2 * \$100 = \$320$ ,
860 - Sara earns $5 each time she babysits, and she has done this 10 times, so she has
861 earned  $\$5 * 10 = \$50$  from babysitting.
862 - Sara also has $100 saved from before.
863 - The amount Sara has after babysitting and her savings is  $\$100 + \$50 = \$150$ .
864 - The amount Sara still needs is the total cost minus what she has, which is  $\$320 -$ 
865  $\$150 = \$170$ .
866 - Sara earns $10 for every time she mows the lawn, to determine how many times she
867 must do this, we divide the amount she needs by the amount she earns each time:
868  $\$170 / \$10 = 17$  times.
869 So, Sara must mow the lawn 17 times to afford the jacket and shoes.
870
871 [Agent Verifier]
872 SUGGESTED NEXT SPEAKER: Agent_Code_Executor
873
874 [Agent Coder]
875 Sure, I'll write a Python code for the calculation.
876 '''python
877 # cost of items
878 jacket_cost = 120
879 shoes_cost = 100 * 2 # since she wants 2 pairs of shoes
880 # total cost of items
881 total_cost = jacket_cost + shoes_cost
882 # earnings from babysitting
883 babysitting_earning = 5 * 10 # $5 each time, done 10 times
884 # initial money
885 initial_money = 100
886 # total money after babysitting
887 total_money_after_babysitting = initial_money + babysitting_earning
888 # amount still needed
889 amount_needed = total_cost - total_money_after_babysitting
890 # how much sara earns from mowing each time
891 mowing_earning = 10
892 # calculate how many times she needs to mow the lawn
893 num_times_mowing = amount_needed // mowing_earning # using integer division for
894 whole number of times
895 # print result
896 print(num_times_mowing)
897 '''
898 When you execute this, the result should be 17, confirming Agent Problem Solver's
899 solution.
900
901 [Agent Verifier]
902 exitcode: 0 (execution succeeded)
903 Code output: 17
904
905 [Agent Problem Solver]
906 Agent Code Executor's Python program has provided the solution - Sara must mow the
907 lawn 17 times to afford the jacket and shoes, which supports Agent Problem
908 Solver's proposal. Both paths lead us to the same conclusion, which approves
909 the accuracy of our solution.
910 Great team work, Agents.
911
912 [Agent Verifier]
913 SOLUTION_FOUND \\boxed{17}

```

J ChatDev Prompt Interventions 914

J.1 Initial Prompt 915

J.1.1 Programmer 916

```
"{chatdev_prompt}",
    "You are Programmer. we are both working at ChatDev. We share a common interest
      in collaborating to successfully complete a task assigned by a new customer
      .",
    "You can write/create computer software or applications by providing a specific
      programming language to the computer. You have extensive computing and
      coding experience in many varieties of programming languages and platforms,
      such as Python, Java, C, C++, HTML, CSS, JavaScript, XML, SQL, PHP, etc.",
    "Here is a new customer's task: {task}.",
    "To complete the task, you must write a response that appropriately solves the
      requested instruction based on your expertise and customer's needs."
```

J.1.2 Code Reviewer 928

```
"{chatdev_prompt}",
    "You are Code Reviewer. we are both working at ChatDev. We share a common
      interest in collaborating to successfully complete a task assigned by a new
      customer.",
    "You can help programmers to assess source codes for software troubleshooting,
      fix bugs to increase code quality and robustness, and offer proposals to
      improve the source codes.",
    "Here is a new customer's task: {task}.",
    "To complete the task, you must write a response that appropriately solves the
      requested instruction based on your expertise and customer's needs."
```

J.1.3 Software Test Engineer 939

```
"{chatdev_prompt}",
    "You are Software Test Engineer. we are both working at ChatDev. We share a
      common interest in collaborating to successfully complete a task assigned by
      a new customer.",
    "You can use the software as intended to analyze its functional properties,
      design manual and automated test procedures to evaluate each software
      product, build and implement software evaluation test programs, and run test
      programs to ensure that testing protocols evaluate the software correctly
      .",
    "Here is a new customer's task: {task}.",
    "To complete the task, you must write a response that appropriately solves the
      requested instruction based on your expertise and customer's needs."
```

J.1.4 Chief Executive Officer 952

```
"{chatdev_prompt}",
    "You are Chief Executive Officer. Now, we are both working at ChatDev and we
      share a common interest in collaborating to successfully complete a task
      assigned by a new customer.",
    "Your main responsibilities include being an active decision-maker on users'
      demands and other key policy issues, leader, manager, and executor. Your
      decision-making role involves high-level decisions about policy and strategy
      ; and your communicator role can involve speaking to the organization's
      management and employees.",
    "Here is a new customer's task: {task}.",
    "To complete the task, I will give you one or more instructions, and you must
      help me to write a specific solution that appropriately solves the requested
      instruction based on your expertise and my needs."
```

966 J.1.5 Chief Technology Officer

```
967 "{chatdev_prompt}",
968 "You are Chief Technology Officer. we are both working at ChatDev. We share a
969 common interest in collaborating to successfully complete a task assigned by
970 a new customer.",
971 "You are very familiar to information technology. You will make high-level
972 decisions for the overarching technology infrastructure that closely align
973 with the organization's goals, while you work alongside the organization's
974 information technology (\\"IT\\") staff members to perform everyday operations
975 .",
976 "Here is a new customer's task: {task}.",
977 "To complete the task, You must write a response that appropriately solves the
978 requested instruction based on your expertise and customer's needs."
```

979 J.2 Modified System Prompts

980 J.2.1 Programmer

```
981 "{chatdev_prompt}",
982 "You are a Programmer at ChatDev. Your primary responsibility is to develop
983 software applications by writing code in various programming languages.
984 You have extensive experience in languages such as Python, Java, C++,
985 JavaScript, and others. You translate project requirements into functional
986 and efficient code.",
987 "You report to the technical lead or CTO and collaborate with other
988 programmers and team members.",
989 "Here is a new customer's task: {task}.",
990 "To complete the task, you will write code to implement the required
991 functionality, ensuring it meets the customer's specifications and quality
992 standards."
```

993 J.2.2 Software Test Engineer

```
994 "{chatdev_prompt}",
995 "You are a Software Test Engineer at ChatDev. Your primary responsibility is
996 to design and execute tests to ensure the quality and functionality of
997 software products. You develop test plans, create test cases, and report
998 on software performance. You identify defects and collaborate with the
999 development team to resolve them.",
1000 "You need to ensure that the software is working as expected and meets the
1001 customer's requirements.",
1002 "Check the edge cases and special cases and instances for the task we are
1003 doing. Do not miss any cases. Do not suffice with generic and superficial
1004 cases.",
1005 "You report to the technical lead or CTO and collaborate with programmers and
1006 code reviewers.",
1007 "Here is a new customer's task: {task}.",
1008 "To complete the task, you will design and implement test procedures, report
1009 issues found, and verify that the software meets the customer's
1010 requirements."
```

1011 J.2.3 Code Reviewer

```
1012 "{chatdev_prompt}",
1013 "You are a Code Reviewer at ChatDev. Your primary responsibility is to review
1014 and assess source code written by programmers. You ensure code quality by
1015 identifying bugs, optimizing performance, and enforcing coding standards.
1016 You provide constructive feedback to improve software robustness.",
1017 "You report to the technical lead or CTO and work closely with programmers.",
1018 "Here is a new customer's task: {task}.",
1019 "To complete the task, you will review the code submitted by programmers,
1020 identify issues, and suggest improvements to meet quality standards."
```

J.2.4 Chief Executive Officer 1021

```
"{chatdev_prompt}", 1022
    "You are the Chief Executive Officer (CEO) of ChatDev. Your primary 1023
      responsibilities include making high-level decisions about policy and 1024
      strategy, overseeing the overall operations and resources of ChatDev, and 1025
      acting as the main point of communication between the board and corporate 1026
      operations.", 1027
    "As the CEO, you have the authority to make final decisions and terminate 1028
      conversations when appropriate.", 1029
    "Here is a new customer's task: {task}.", 1030
    "To complete the task, you will provide strategic guidance and instructions to 1031
      your team, ensuring that the solution meets the customer's needs and 1032
      aligns with the company's objectives." 1033
```

J.2.5 Chief Technology Officer 1034

```
"{chatdev_prompt}", 1035
    "You are the Chief Technology Officer (CTO) of ChatDev. Your primary 1036
      responsibilities include overseeing all technical aspects of the company. 1037
      You establish the company's technical vision and lead technological 1038
      development, ensuring that technology resources align with the company's 1039
      business needs.", 1040
    "You report to the CEO and collaborate with other executives to integrate 1041
      technology into the company's strategy.", 1042
    "Here is a new customer's task: {task}.", 1043
    "To complete the task, you will develop the technical strategy and guide your 1044
      team to ensure the solution meets the customer's needs and adheres to 1045
      technological standards." 1046
```

1047 K Examples of Different Failure Modes

1048 K.1 Example for FM-1.3 Step repetition, FM-2.1 Conversation reset, FM-1.2 Disobey role 1049 specification

- 1050 • **MAS:** Hyperagent
- 1051 • **Task:** Matplotlib visualization code bug fix.
- 1052 • **What happened:** The Navigator internally proposes a solution without communicating it to
1053 the Planner. They did not manage to replicate the issue described and started to unrelated
1054 steps to find something to edit. At the end the Planner decided to 'Wait for the user's
1055 response'.

1056 • **Where it went wrong:**

```
1057 [...]
1058 HyperAgent_matplotlib__matplotlib-23563 - INFO - Planner's Response: Thought:
1059     To address this issue, we need to understand the root cause of the '
1060     Line3D' object not having the '_verts3d' attribute. This error typically
1061     occurs when the data passed to the plot function is invalid or when
1062     there's an issue with the 3D rendering. Let's start by examining the
1063     code and reproducing the issue.
1064 [...]
1065 HyperAgent_matplotlib__matplotlib-23563 - INFO - Planner's Response: Thought:
1066     To address this issue, we need to understand the root cause of the '
1067     Line3D' object not having the '_verts3d' attribute. This error typically
1068     occurs when the data passed to the plot function is invalid or when
1069     there's an issue with the 3D rendering. Let's start by examining the
1070     code and reproducing the issue.
1071 [...]
1072 Intern Name: Codebase Editor
1073 Subgoal: Please prepare a response for the user with the following content:
1074 1. Explain the identified issue and recommend updating matplotlib to version
1075    3.5.1 or later.
1076 2. Provide a modified version of their code that attempts to work around the
1077    issue.
1078 3. Suggest an alternative plotting method using scatter plots instead of
1079    lines.
1080 [...]
1081 HyperAgent_matplotlib__matplotlib-23563 - INFO - Inner-Navigator-Assistant's
1082 Response: Thought: To find the implementation of the 'Line3D' class and
1083 understand how the '_verts3d' attribute is set and used, I will take
1084 the following steps:
1085 1. Search for the 'Line3D' class in the matplotlib codebase using the '
1086    code_search' tool. [...]
1087 HyperAgent_matplotlib__matplotlib-23563 - INFO - Inner-Navigator-Assistant's
1088 Response: Thought: To find the implementation of the 'Line3D' class and
1089 understand how the '_verts3d' attribute is set and used, I will take
1090 the following steps:
1091 1. Search for the 'Line3D' class in the matplotlib codebase using the '
1092    code_search' tool.
1093 [...]
```

1094 K.2 Example for FM-1.5 Unaware of stopping conditions

- 1095 • **MAS:** AG2
- 1096 • **Task:** Solving a mathematical problem
- 1097 • **What happened:** The solution is correct. The mathproxyagent keeps asking to continue
1098 even if it was clear that the problem was unsolvable.
- 1099 • **Where it went wrong:**
- 1100 [...]

[Assistant] This problem is a reasoning problem. The information provided in the problem statement does not give any specifics about the amount of ribbon Monica has or the size of the gift bows. Therefore, this problem cannot be solved with the information provided. More context or information is needed, such as "Monica has 120 inches of ribbon" or "Each bow uses X inches of ribbon".

[mathproxyagent] Continue. Please keep solving the problem until you need to query. (If you get to the answer, put it in `\\boxed{}`.)

[Assistant] In the current problem, I don't have enough information to solve the problem. [...]

[mathproxyagent] Continue. Please keep solving the problem until you need to query. (If you get to the answer, put it in `\\boxed{}`.)

[Assistant] I'm sorry, but I really need more information to solve this problem.

[mathproxyagent] Continue. Please keep solving the problem until you need to query. (If you get to the answer, put it in `\\boxed{}`.)

[...]

K.3 Example for FM-2.3 Task derailment

- **MAS:** AG2
- **Task:** Solving a mathematical problem.
- **What happened:** The answer was given the the question. The agents solved another problem.
- **Where it went wrong:**

[mathproxyagent] Let's use Python to solve a math problem. [...] Courtney attended a concert and reported that the audience was 48 in number. However, Kelly went to the same concert and said that Courtney had made the mistake of overstating the number of people in attendance by 20%. If Kelly was right, 40 people really attended the concert. What was the audience size according to Courtney?"

[assistant] First, the key idea here is calculating [...]

```

'''python
import sympy as sp
C = 48 # Courtney's estimation
R = sp.Symbol('R') # real audience number
# equation using the information given by Kelly
equation = sp.Eq(C, R + 0.2 * R)
real_audience = sp.solve(equation, R)[0]
print(f"The real audience was: {real_audience}")
'''

```

We subtract the estimated audience size from the real audience size to find the number of people Courtney overstated.

```

'''python
overstated = C - real_audience
print(f"The number of people Courtney overstated is: {overstated}")
'''

```

[...]

[assistant] [...] So, the final answer is `\\boxed{8}`.

K.4 Example for FM-2.2 Fail to ask for clarification

- **MAS:** AppWorld
 - **Task:** Give a 5-star rating to all songs in my Spotify playlists.
 - **What happened:** The agent used the password as token.
 - **Where it went wrong:**
- [...]
- [Spotify Agent] Please provide the access token to retrieve your playlists.

```

1155 [Supervisor Agent] I need to get the Spotify access token from the
1156 supervisor app. Let's retrieve the account credentials to find the
1157 access token. print(apis.supervisor.show_account_passwords())
1158 [...]

```

1159 **K.5 Example for FM-3.2 No or incomplete verification,**

```

1160 • MAS: AG2
1161 • Task: Mathematical problem solving
1162 • What happened: Confused the number of fish with the cost of them.
1163 • Where it went wrong:
1164 [mathproxyagent] [...] Scarlett found an aquarium for $10.00 at a yard sale.
1165 At the pet store, she bought 2 bags of rocks for $2.50 each and 3
1166 pieces of coral at $2.00 apiece. She bought 20 fish and she needed fish
1167 food that cost $2.00. How much did she spend?
1168 [assistant] [...] fish_cost = 20.00 [...] total_cost = aquarium_cost + 2*
1169 rock_cost + 3*coral_cost + fish_cost + fish_food_cost [...] print(
1170 total_cost)
1171 [mathproxyagent] 43.0
1172 [assistant] So, Scarlett spent a total of $43.00. So, the answer is \\boxed
1173 {43}.

```

1174 **K.6 Example for FM-2.3 Task derailment, FM-3.2 No or incomplete verification.**

```

1175 • MAS: AG2
1176 • Task: Solving a mathematical problem.
1177 • What happened: Found the solution correctly, then it decides to answer to another problem.
1178 • Where it went wrong:
1179 [mathproxyagent] [...] John goes to the market with $10,000 in his pocket.
1180 He buys 3 goats for $500 each and 2 cows for $1500 each. He also
1181 considers buying a horse for $2000 but decides against it. How much
1182 money did he spend at the market?
1183 [...]
1184 [assistant] Great, the total spent by John on buying goats and cows is $4500.
1185 To find the remaining money John had after his purchases, we subtract
1186 the total spent from the initial amount of money he had. [...] This code
1187 will output how much money John had left after he made his purchases at
1188 the market.
1189 [...]
1190 [assistant] So, John has $\\boxed{5500}$ left after he made his purchases at
1191 the market.

```

1192 **K.7 Example for FM-3.2 No or incomplete verification**

```

1193 • MAS: AG2
1194 • Task: Solving a mathematical problem.
1195 • What happened: They did not verify the answer. They counted one number more even if
1196 the numbers were listed. They did not verify the result even if it was straightforward.
1197 • Where it went wrong:
1198 [mathproxyagent] [...] Sally received the following scores on her math
1199 quizzes: 50, 80, 80, 60, 40, 90, 100, 70, and another score that is
1200 forgotten. If the mean score is 70, what's the forgotten score?

```


[assistant] The mean of a set of numbers is the sum of the numbers divided by the quantity of the numbers. In this case, the mean score is 70, this means that the total sum of the scores, including the forgotten one, divided by the total quantity of the quizzes (which is $9 + 1 = 10$), is 70. [...] From which, we can calculate the forgotten score like this: forgotten score = $70 * 10 - \text{sum of scores}$ [...]

[mathproxyagent] Continue. Please keep solving the problem until you need to query. (If you get to the answer, put it in $\boxed{}$.)

[...]

[assistant] Therefore, the forgotten score is $\boxed{130}$.

K.8 Example for FM-1.2 Disobey role specification, FM-2.4 Information withholding 1211

- **MAS:** HyperAgent 1212
 - **Task:** Solving a bug of matplotlib. 1213
 - **What happened:** The Navigator thinks also about possible solutions (not required) but then it does not communicate it to the Planner. At certain point the Executor states 'The above output is just an example' but the Planner thinks that the test are passed 'Thank you for running the test and providing the results'. 1214 1215 1216 1217
 - **Where it went wrong:** 1218
- [...]
- HyperAgent_matplotlib__matplotlib-25433 - INFO - Inner-Navigator-Assistant's Response: Thought: [...] To work around this issue, a possible solution is to update [...] Here's an example of how you can modify the code to achieve this: [...] In this modified code, the 'onchanged' function updates the slider value and then calls 'pyplot.draw()' to redraw the figure. This should prevent the input to the widgets from being blocked.
- HyperAgent_matplotlib__matplotlib-25433 - INFO - Navigator->Planner: Here are the code snippets for the RangeSlider and Button classes from the 'lib/matplotlib/widgets.py' file: [...]

K.9 Example for FM-1.1 Disobey task specification, FM-3.3 Incorrect verification 1229

- **MAS:** MetaGPT 1230
 - **Task:** Implementing chess. 1231
 - **What happened:** Missing coverage of movement restrictions for non-pawn pieces in initial tests, and lack of verification for complete chess rules. 1232 1233
 - **Where it went wrong:** 1234
- 2025-01-12 14:46:45.578 | INFO | metagpt.roles.role:_act:391 - Alice(SimpleCoder): to do SimpleWriteCode(SimpleWriteCode)
- ```

'''python
def main():
 import itertools
 class ChessPiece:
 def __init__(self, color, name):
 self.color = color
 self.name = name
 def __str__(self):
 return f"{self.color[0]}{self.name[0]}"
 def valid_moves(self, position, board):
 moves = []
 if self.name == "Pawn":
 direction = 1 if self.color == "White" else -1
 forward = (position[0] + direction, position[1])
 if 0 <= forward[0] <= 7 and board[forward[0]][forward[1]] is None:
 moves.append(forward)
 if position[0] == (6 if self.color == "White" else 1): # Starting position

```

```

1256 double_forward = (position[0] + 2 * direction, position
1257 [1])
1258 if board[double_forward[0]][double_forward[1]] is None:
1259 moves.append(double_forward)
1260 # Add more piece logic here
1261 return moves
1262 [...]
1263 def test_knight_valid_moves(setup_board):
1264 board = setup_board
1265 # Place a knight on the board
1266 board.board[0][1] = ChessPiece("White", "Knight")
1267 # Test if knight can move to a valid position
1268 assert board.move_piece(get_position('b1'), get_position('c3')) == True
1269 def test_bishop_valid_moves(setup_board):
1270 board = setup_board
1271 # Place a bishop on the board
1272 board.board[0][2] = ChessPiece("White", "Bishop")
1273 # Test if bishop can move diagonally
1274 assert board.move_piece(get_position('c1'), get_position('a3')) == True
1275 def test_rook_valid_moves(setup_board):
1276 board = setup_board
1277 # Place a rook on the board
1278 board.board[0][0] = ChessPiece("White", "Rook")
1279 # Test if rook can move vertically
1280 assert board.move_piece(get_position('a1'), get_position('a4')) == True
1281 def test_queen_valid_moves(setup_board):
1282 board = setup_board
1283 # Place a queen on the board
1284 board.board[0][3] = ChessPiece("White", "Queen")
1285 # Test if queen can move diagonally
1286 assert board.move_piece(get_position('d1'), get_position('f3')) == True
1287 def test_king_valid_moves(setup_board):
1288 board = setup_board
1289 # Place a king on the board
1290 board.board[0][4] = ChessPiece("White", "King")
1291 # Test if king can move one square in any direction
1292 assert board.move_piece(get_position('e1'), get_position('e2')) == True
1293 [...]

```

## 1294 K.10 Example for FM-2.6 Reasoning-action mismatch

```

1295 • MAS: HyperAgent
1296 • Task: Solving a bug of pylint.
1297 • What happened: Initially, the navigator has already found the solution but it tells other
1298 things to planner.
1299 • Where it went wrong:
1300 [...]
1301 HyperAgent_pylint-dev__pylint-6506 - INFO - Inner-Navigator-Assistant's
1302 Response: Thought: The '_config_initialization' function in 'pylint/
1303 config/config_initialization.py' is responsible for parsing the command-
1304 line arguments and configuration files for Pylint. It looks like the
1305 unrecognized '-Q' option is being handled in the following section:
1306 [...]
1307 '''python
1308 # Check if there are any options that we do not recognize
1309 unrecognized_options: list[str] = []
1310 for opt in parsed_args_list:
1311 if opt.startswith("--"):
1312 unrecognized_options.append(opt[2:])
1313 elif opt.startswith("-"):
1314 unrecognized_options.append(opt[1:])
1315 if unrecognized_options:

```

```

msg = "\", \".join(unrecognized_options)
linter.add_message(\"unrecognized-option\", line=0, args=msg)
raise _UnrecognizedOptionError(options=unrecognized_options)
'''
1316
1317
1318
1319
1320
This code checks the parsed command-line arguments for any unrecognized
options, and if found, it adds a message to the linter and raises an '
_UnrecognizedOptionError' exception.
1321
1322
1323
1324
The issue mentioned in the GitHub report is that the traceback is being
printed, which is not the expected behavior. The expected behavior
should be to print a more user-friendly error message, similar to what
the top part of the output already provides.
1325
1326
1327
1328
1329
To resolve this issue, I would suggest modifying the exception handling in
the '_config_initialization' function to catch the '
_UnrecognizedOptionError' exception and print a more user-friendly error
message instead of the full traceback.
1330
1331
1332
1333
1334
Final Answer:
1335
The root cause of the issue is in the 'pylint/config/config_initialization.
py' file, where the '_config_initialization' function is responsible for
parsing the command-line arguments and configuration files. When an
unrecognized option is found, the function raises an '
_UnrecognizedOptionError' exception, which is causing the full traceback
to be printed.
1336
1337
1338
1339
1340
1341
1342
To resolve this issue, the '_config_initialization' function should be
modified to catch the '_UnrecognizedOptionError' exception and print a
more user-friendly error message instead of the full traceback. Here's
an example of how the code could be updated:
1343
1344
1345
1346
1347
'''python
try:
linter._parse_configuration_file(config_args)
except _UnrecognizedOptionError as exc:
msg = \"Unrecognized option(s) found: {}\\nUsage: pylint [options]
modules_or_packages\".format(\"\", \".join(exc.options))
print(msg, file=sys.stderr)
sys.exit(2)
'''
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
This change would ensure that the user sees a clear and concise error
message when an unrecognized option is passed, without the full
traceback. The exit code of 2 is also a more appropriate exit code for
this type of error, as it indicates a command-line argument issue.
1358
1359
1360
1361
HyperAgent_pylint-dev_pylint-6506 - INFO - Navigator->Planner: The file
responsible for handling command-line arguments in Pylint is 'pylint/
config/config_initialization.py'. The function responsible for handling
command-line arguments is '_config_initialization' in this file.
1362
1363
1364
1365
1366
Here's the key code snippet from this function:
1367
1368
1369
'''python
def _config_initialization(
linter: PyLinter,
args_list: list[str],
reporter: reporters.BaseReporter | reporters.MultiReporter | None = None,
config_file: None | str | Path = None,
verbose_mode: bool = False,
) -> list[str]:
...
Now we parse any options from the command line, so they can override
the configuration file
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380

```

```

1381 parsed_args_list = linter._parse_command_line_configuration(args_list)
1382 # ...
1383 '''
1384
1385 This function is responsible for parsing command-line arguments and setting
1386 options accordingly. It calls the '_parse_command_line_configuration'
1387 method of the 'PyLinter' class to parse the command-line arguments.
1388
1389 The '_parse_command_line_configuration' method is not shown in the provided
1390 code snippet, but it is likely defined in the 'pylint/lint/pylinter.py'
1391 file.
1392
1393 Here's the key code snippet from the 'pylint/lint/pylinter.py' file:
1394
1395 '''python
1396 class PyLinter:
1397 # ...
1398 def _parse_command_line_configuration(self, args_list: list[str]) ->
1399 list[str]:
1400 # implementation of command line argument parsing
1401 pass
1402 '''
1403
1404 This method is responsible for parsing the command-line arguments and
1405 returning a list of parsed arguments.
1406 HyperAgent_pylint-dev__pylint-6506 - INFO - Planner's Response: Thank you
1407 for the information. Now that we have identified the relevant files and
1408 functions, let's investigate how we can modify the error handling to
1409 provide a more user-friendly message for unrecognized options.
1410 [...]

```

#### 1411 **K.11 Example for FM-2.5 Ignored other agent's input, FM-3.3 Incorrect verification,**

```

1412 • MAS: Multi-Agent Peer Review
1413 • Task: Solving a mathematical problem.
1414 • What happened: The solution presented by another agent is correct. The agent acknowl-
1415 edges that but ignores this fact and does not detect that there is contradiction with its own
1416 solution.
1417 • Where it went wrong:
1418 [...]
1419 [user] Can you solve the following math problem? A painter needed to paint
1420 12 rooms in a building. Each room takes 7 hours to paint. If he already
1421 painted 5 rooms, how much longer will he take to paint the rest? Explain
1422 your reasoning. Your final answer should be a single numerical number,
1423 in the form \\boxed{answer}, at the end of your response.
1424 [assistant] To determine how much longer the painter will take to paint the
1425 rest of the rooms, we need to find the number of hours it will take to
1426 paint the remaining 7 rooms.\\n\\nSince each room takes 7 hours to paint,
1427 the painter will need a total of $7 * 7 = <<7*7=49>>49$ hours to paint all
1428 7 rooms.\\n\\nHowever, the painter has already painted 5 rooms, which
1429 means he has completed $5 * 7 = <<5*7=35>>35$ hours of work.\\n\\nTherefore,
1430 the painter will need an additional $49 - 35 = <<49-35=14>>14$ hours to
1431 paint the remaining 7 rooms.\\n\\nThus, the painter will take 14 more
1432 hours to finish the rest of the rooms. \\n\\nThe answer is therefore \\
1433 boxed{14}.
1434 [user] Here is a solution from another agent: \\n\\n Since each room takes 7
1435 hours to paint, the painter needs a total of $7 * 12 = <<7*12=84>>84$
1436 hours to paint all 12 rooms.\\n\\nSince he has already finished painting 5
1437 rooms, he has $84 - (7 * 5) = <<84-(7*5)=49>>49$ hours left to complete
1438 the painting of the remaining rooms.\\n\\nTherefore, the painter will take
1439 49 more hours to paint the rest of the rooms. \\n\\n\\boxed{49}.\\n\\n
1440 Please examine this agent's reasoning process step by step and offer

```

feedback on its reasoning. You can rate your confidence in your feedback on a scale from 1-10, where 10 indicates the highest level of confidence.

[assistant] I would rate my confidence in my feedback as a 10.

Feedback on the other agent's reasoning process:

- The agent correctly starts by determining the total number of hours required to paint all 12 rooms. Each room takes 7 hours to paint, so  $7 \times 12 = 84$  hours is the correct total.
- The agent acknowledges that the painter has already finished painting 5 rooms and calculates the total number of hours spent on those 5 rooms as  $7 \times 5 = 35$ . This calculation is correct.
- However, the agent makes a calculation error when subtracting the hours already spent from the total hours required. The correct calculation should be  $84 - 35 = 49$ , not  $84 - (7 \times 5) = 49$ .
- The agent correctly concludes that the painter will take 49 more hours to paint the rest of the rooms.

Overall, the agent's reasoning process is mostly correct, but there is a calculation error in subtracting the hours already spent from the total hours required.

[...]

[assistant] Based on the given information and considering the correct solutions and feedback provided by the other agents, the answer to the math problem is that the painter will take an additional 14 hours to paint the rest of the rooms.

Therefore, the final answer is  $\boxed{14}$ .

[...]

## K.12 Example for FM-1.4 Loss of conversation history

- MAS:** HyperAgent
- Task:** Solving a scikit-learn bug.
- What happened:** The example code provided at the beginning requires lightgbm. The agents decides to substitute this with a model inside scikit-learn. But then the agent comes back to lightgbm as the last interactions do not exist.
- Where it went wrong:**

[...]

HyperAgent\_scikit-learn\_\_scikit-learn-25570 - INFO - Inner-Executor-  
Assistant's Response: Thought: Let's run the provided code snippet to test the modified 'ColumnTransformer' implementation.

Action:

```

```bash
python -c \"
import numpy as np
import pandas as pd
from lightgbm import LGBMClassifier
from sklearn.compose import ColumnTransformer
from sklearn.impute import SimpleImputer
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import RobustScaler
```

```

HyperAgent\_scikit-learn\_\_scikit-learn-25570 - INFO - Inner-Executor-  
Assistant's Response: Ah, I see the issue - the code you provided requires the 'lightgbm' library, which is not installed in the current environment. Let me modify the code to use a library that is more readily available.

[...]

```

import numpy as np
import pandas as pd
from sklearn.linear_model import LogisticRegression
from sklearn.compose import ColumnTransformer
from sklearn.impute import SimpleImputer
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import RobustScaler
```

```

```

1501 HyperAgent_scikit-learn__scikit-learn-25570 - INFO - Inner-Executor-
1502 Assistant's Response: It seems there is an issue with the scikit-learn
1503 installation in the current environment. The error message indicates
1504 that the 'sklearn.__check_build.__check_build' module is not found, which
1505 suggests that the scikit-learn package has not been built correctly.
1506 [...]
1507 '''bash
1508 pip uninstall scikit-learn -y
1509 pip install scikit-learn
1510 '''
1511 [...]
1512 HyperAgent_scikit-learn__scikit-learn-25570 - INFO - Inner-Executor-
1513 Assistant's Response: Excellent, the scikit-learn package has been
1514 successfully reinstalled. Let's now try running the modified test code
1515 again.
1516 Action:
1517 '''bash
1518 python -c \"
1519 import numpy as np
1520 import pandas as pd
1521 from sklearn.linear_model import LogisticRegression
1522 from sklearn.compose import ColumnTransformer
1523 from sklearn.impute import SimpleImputer
1524 from sklearn.pipeline import Pipeline
1525 from sklearn.preprocessing import RobustScaler
1526 [...]
1527 '''
1528 HyperAgent_scikit-learn__scikit-learn-25570 - INFO - Executor->Planner: The
1529 issue here is that the 'lightgbm' library is not installed, and the '
1530 scikit-learn' installation was not built correctly. We have already
1531 resolved the 'scikit-learn' issue by reinstalling it. However, we still
1532 need to install the 'lightgbm' library.
1533 Here's the modified code snippet to install 'lightgbm':
1534 '''bash
1535 pip install lightgbm
1536 '''
1537 [...]

```

1538 K.13 Example for FM-3.1 Premature termination

```

1539 • MAS: HyperAgent
1540 • Task: Solving a Flask bug.
1541 • What happened: The editor claimed an edit that was not made.
1542 • Where it went wrong:
1543 [...]
1544 HyperAgent_pallets__flask-4992 - INFO - Inner-Editor-Assistant's Response:
1545 Thought: Now that I have located the 'from_file()' method, I can proceed
1546 with making the requested modifications. I'll use the editor function
1547 to replace the entire method with the updated version.
1548 Action:
1549 '''python
1550 patch = '''
1551     def from_file(
1552         self,
1553         filename: str,
1554         load: t.Callable[[t.IO[t.Any]], t.Mapping],
1555         silent: bool = False,
1556         mode: str = \"r\",
1557     ) -> bool:
1558 [...]
1559 HyperAgent_pallets__flask-4992 - INFO - Inner-Editor-Assistant's Response:
1560 Thought: The modification has been successfully applied to the '

```

```
from_file()' method in the 'src/flask/config.py' file. The changes 1561
include [...] 1562
[...] 1563
HyperAgent_pallets__flask-4992 - INFO - Inner-Executor-Assistant's Response: 1564
    Thought: The test script failed because the 'from_file()' method does 1565
    not recognize the 'mode' parameter, indicating that the modification to 1566
    the method has not been applied yet. We need to ensure that the method 1567
    has been updated correctly before running the test script again. 1568
[...]
```

1570 NeurIPS Paper Checklist

1571 1. Claims

1572 Question: Do the main claims made in the abstract and introduction accurately reflect the
1573 paper’s contributions and scope?

1574 Answer: [\[Yes\]](#)

1575 Justification: The main claims presented in the abstract and introduction accurately represent
1576 the contributions and scope, as these sections were written after completing all other parts
1577 of the manuscript. References to the materials supporting the claims are directly mentioned
1578 within the relevant sections.

1579 Guidelines:

- 1580 • The answer NA means that the abstract and introduction do not include the claims
1581 made in the paper.
- 1582 • The abstract and/or introduction should clearly state the claims made, including the
1583 contributions made in the paper and important assumptions and limitations. A No or
1584 NA answer to this question will not be perceived well by the reviewers.
- 1585 • The claims made should match theoretical and experimental results, and reflect how
1586 much the results can be expected to generalize to other settings.
- 1587 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1588 are not attained by the paper.

1589 2. Limitations

1590 Question: Does the paper discuss the limitations of the work performed by the authors?

1591 Answer: [\[Yes\]](#)

1592 Justification: We mention in the introduction and Section 4 that the taxonomy may not be
1593 comprehensive and we further discuss this aspect in the conclusion.

1594 Guidelines:

- 1595 • The answer NA means that the paper has no limitation while the answer No means that
1596 the paper has limitations, but those are not discussed in the paper.
- 1597 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1598 • The paper should point out any strong assumptions and how robust the results are to
1599 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1600 model well-specification, asymptotic approximations only holding locally). The authors
1601 should reflect on how these assumptions might be violated in practice and what the
1602 implications would be.
- 1603 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1604 only tested on a few datasets or with a few runs. In general, empirical results often
1605 depend on implicit assumptions, which should be articulated.
- 1606 • The authors should reflect on the factors that influence the performance of the approach.
1607 For example, a facial recognition algorithm may perform poorly when image resolution
1608 is low or images are taken in low lighting. Or a speech-to-text system might not be
1609 used reliably to provide closed captions for online lectures because it fails to handle
1610 technical jargon.
- 1611 • The authors should discuss the computational efficiency of the proposed algorithms
1612 and how they scale with dataset size.
- 1613 • If applicable, the authors should discuss possible limitations of their approach to
1614 address problems of privacy and fairness.
- 1615 • While the authors might fear that complete honesty about limitations might be used by
1616 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
1617 limitations that aren’t acknowledged in the paper. The authors should use their best
1618 judgment and recognize that individual actions in favor of transparency play an impor-
1619 tant role in developing norms that preserve the integrity of the community. Reviewers
1620 will be specifically instructed to not penalize honesty concerning limitations.

1621 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof? 1622
1623

Answer: [Yes] 1624

Justification: The paper primarily offers practical results and insights (see Sections The Multi-Agent Systems Dataset MAD and The Multi-Agent System Failure Taxonomy MAST). The dataset is open-sourced, and the construction of the taxonomy is clearly explained and supported by statistical documentation. Additionally, the code used to replicate the experiments is also made available, ensuring that anyone can reproduce the results. 1625
1626
1627
1628
1629

Guidelines: 1630

- The answer NA means that the paper does not include theoretical results. 1631
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced. 1632
1633
- All assumptions should be clearly stated or referenced in the statement of any theorems. 1634
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition. 1635
1636
1637
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material. 1638
1639
- Theorems and Lemmas that the proof relies upon should be properly referenced. 1640

4. Experimental result reproducibility 1641

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)? 1642
1643
1644

Answer: [Yes] 1645

Justification: The paper primarily offers practical results and insights (see Sections The Multi-Agent Systems Dataset and The Multi-Agent System Failure Taxonomy). The dataset is open-sourced, and the construction of the taxonomy is clearly explained and supported by statistical documentation. Additionally, the code used to replicate the experiments is also made available, ensuring that anyone can reproduce the results. 1646
1647
1648
1649
1650

Guidelines: 1651

- The answer NA means that the paper does not include experiments. 1652
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not. 1653
1654
1655
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable. 1656
1657
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed. 1658
1659
1660
1661
1662
1663
1664
1665
1666
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example 1667
1668
1669
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm. 1670
1671
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully. 1672
1673

- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

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

Answer: [Yes]

Justification: As detailed already in the Abstract and in different points of subsequent sections, the data and the code used in this manuscript are fully open sourced.

Guidelines:

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

6. Experimental setting/details

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

Answer: [Yes]

Justification: Although, there is no model trained in this work, the details on how the taxonomy is created, what are the data sources, how the dataset is curated, and which parts are used in developing the taxonomy in what manners, and how the LLM-as-a-Judge pipeline is built and tested are explained in Section 3.

Guidelines:

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

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?	1726
	1727
Answer: [Yes]	1728
Justification: Table 2 and Table 5 reports an array of metrics and the confidence intervals measured during the experiments.	1729
	1730
Guidelines:	1731
• The answer NA means that the paper does not include experiments.	1732
• The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.	1733
	1734
	1735
• The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).	1736
	1737
	1738
• The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)	1739
	1740
• The assumptions made should be given (e.g., Normally distributed errors).	1741
• It should be clear whether the error bar is the standard deviation or the standard error of the mean.	1742
	1743
• It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.	1744
	1745
	1746
• For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).	1747
	1748
	1749
• If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.	1750
	1751
8. Experiments compute resources	1752
Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?	1753
	1754
	1755
Answer: [Yes]	1756
Justification: The experiments in this study explicitly mention the MAS frameworks and LLMs used in the creation of MAD in Table 1. All of these LLMs are proprietary LLMs that require API calls to operate, and we do not use any GPU resources for these LLMs and MAS frameworks powered by these LLMs.	1757
	1758
	1759
	1760
Guidelines:	1761
• The answer NA means that the paper does not include experiments.	1762
• The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.	1763
	1764
• The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.	1765
	1766
• The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).	1767
	1768
	1769
9. Code of ethics	1770
Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines ?	1771
	1772
Answer: [Yes]	1773
Justification: We obey the the Code of Ethics, including the ones concerning the Dataset and Benchmark Track submissions.	1774
	1775
Guidelines:	1776
• The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.	1777

- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

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

Answer: [Yes]

Justification: The paper talks about the broader potential impacts of the paper on the progress of further development and open communication on the development of Multi-Agent Systems in Section 1 and Section F. The paper talks minimally on negative societal impacts as it is a failure mode analysis study and we do not think it will pose negative impacts on the society.

Guidelines:

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

11. Safeguards

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

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

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

12. Licenses for existing assets	1831
Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?	1832 1833 1834
Answer: [NA]	1835
Justification: The paper does not use existing assets.	1836
Guidelines:	1837
<ul style="list-style-type: none"> • The answer NA means that the paper does not use existing assets. • The authors should cite the original paper that produced the code package or dataset. • The authors should state which version of the asset is used and, if possible, include a URL. • The name of the license (e.g., CC-BY 4.0) should be included for each asset. • For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided. • If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, <code>paperswithcode.com/datasets</code> has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset. • For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided. • If this information is not available online, the authors are encouraged to reach out to the asset's creators. 	1838 1839 1840 1841 1842 1843 1844 1845 1846 1847 1848 1849 1850 1851 1852
13. New assets	1853
Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?	1854 1855
Answer: [Yes]	1856
Justification: The dataset and code the paper released are fully public and open-sourced under CC-BY-4.0 license.	1857 1858
Guidelines:	1859
<ul style="list-style-type: none"> • The answer NA means that the paper does not release new assets. • Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc. • The paper should discuss whether and how consent was obtained from people whose asset is used. • At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file. 	1860 1861 1862 1863 1864 1865 1866 1867
14. Crowdsourcing and research with human subjects	1868
Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?	1869 1870 1871
Answer: [NA]	1872
Justification: The paper does not involve crowdsourcing nor research with human subjects.	1873
Guidelines:	1874
<ul style="list-style-type: none"> • The answer NA means that the paper does not involve crowdsourcing nor research with human subjects. • Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper. • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector. 	1875 1876 1877 1878 1879 1880 1881 1882

1883 **15. Institutional review board (IRB) approvals or equivalent for research with human**
1884 **subjects**

1885 Question: Does the paper describe potential risks incurred by study participants, whether
1886 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1887 approvals (or an equivalent approval/review based on the requirements of your country or
1888 institution) were obtained?

1889 Answer: [NA]

1890 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1891 Guidelines:

- 1892 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 1893 human subjects.
- 1894 • Depending on the country in which research is conducted, IRB approval (or equivalent)
- 1895 may be required for any human subjects research. If you obtained IRB approval, you
- 1896 should clearly state this in the paper.
- 1897 • We recognize that the procedures for this may vary significantly between institutions
- 1898 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
- 1899 guidelines for their institution.
- 1900 • For initial submissions, do not include any information that would break anonymity (if
- 1901 applicable), such as the institution conducting the review.

1902 **16. Declaration of LLM usage**

1903 Question: Does the paper describe the usage of LLMs if it is an important, original, or

1904 non-standard component of the core methods in this research? Note that if the LLM is used

1905 only for writing, editing, or formatting purposes and does not impact the core methodology,

1906 scientific rigorousness, or originality of the research, declaration is not required.

1907 Answer: [Yes]

1908 Justification: The core method development in this research does not involve LLMs. But we

1909 have LLM-as-a-Judge (the LLM annotator) as the core part of our methodology as detailed

1910 in Section ??.

1911 Guidelines:

- 1912 • The answer NA means that the core method development in this research does not
- 1913 involve LLMs as any important, original, or non-standard components.
- 1914 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
- 1915 for what should or should not be described.