

Evaluation and Benchmarking of Multi-Agent LLM Systems: A Comprehensive Review

Yash Agrawal

yash.agr96@gmail.com

Abstract:

We are seeing exciting new possibilities as systems built on large language models (LLMs) begin to work together in teams, collaborating, negotiating, and coordinating to tackle complex tasks. However, figuring out how to properly evaluate these multi-agent systems is still a work in progress. Many current approaches rely on overly simple tasks, scattered observations, or domain-specific benchmarks. This review takes a closer look at how these systems are currently being evaluated, introduces a framework for understanding different evaluation needs, and highlights major gaps. It also suggests a path forward toward more consistent and reliable benchmarking. Clear, structured evaluation methods are essential for measuring how well these systems collaborate, adapt, scale, and interact with humans. Establishing shared standards will help advance the field, make results easier to compare, and support the safe and effective use of multi-agent systems.

Keywords: Multi-agent systems, large language models, emergent behavior, robustness and safety, human–AI teaming, scalability and efficiency, collective intelligence, trust and interpretability, standardized metrics, synthetic societies, cross-domain evaluation, reproducibility.

1. INTRODUCTION

Multi-agent systems have been explored for decades in fields like robotics, distributed systems, and game theory. Traditionally, these systems involved agents that were highly specialized and followed fixed communication rules. But with the rise of large language models (LLMs), the landscape has changed. Agents can now communicate and coordinate using natural language, opening the door to more flexible, dynamic collaboration.

This shift brings up an important question: how should we evaluate these new types of multi-agent systems?

For single-agent LLMs (e.g., MMLU [1], HELM [2], BIG-Bench [3]), progress has traditionally been evaluated using standardized benchmarks. In contrast, evaluation methods for multi-agent systems are still scattered, ranging from small simulations to domain-specific case studies and informal observations of emergent behavior.

In this review, we take a closer look at current evaluation methods for multi-agent systems, introduce a framework for understanding different evaluation dimensions, highlight key gaps in existing approaches, and suggest a roadmap for developing more consistent and meaningful benchmarks for multi-agent LLM systems.

2. BACKGROUND: FROM CLASSICAL MAS TO LLM-MAS

2.1 Classical MAS Evaluation

In its earlier forms, research on multi-agent systems (MAS) focused on clearly defined tasks and environments. For Example **Game-theoretic models** such as the Prisoner’s Dilemma or Nash bargaining were widely used to study strategic interaction, cooperation, and competition between agents. With the rise of reinforcement learning, benchmarks expanded to include richer testbeds like the **StarCraft Multi-**

Agent Challenge[4] and platforms such as **PettingZoo**[5], which provided more dynamic settings for coordination and control.

In parallel, **robotics and swarm intelligence** offered physical demonstrations of MAS performance. Here, evaluation typically relied on tangible measures such as spatial coverage, coordination efficiency, or the ability of a swarm to reach consensus. While effective within their domains, these approaches were designed for structured problems and provided limited insight into more complex or unpredictable behaviors.

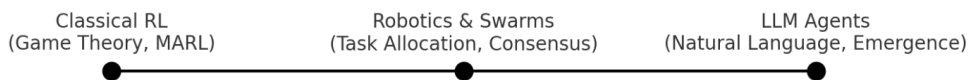
2.2 LLM-Enabled Multi-Agent Systems

Large language models have shifted the ground beneath MAS research. Unlike earlier systems that depended on rigid communication protocols, LLM-based agents can coordinate using **natural language**. This makes interaction more flexible and human-like, but also introduces ambiguity and challenges in interpretation.

The range of tasks has also broadened. Instead of focusing only on navigation or resource-sharing, agents can now take on **open-ended problems** such as multi-party negotiation, collaborative planning, or even joint storytelling. Perhaps most intriguing is the rise of emergent behaviors, ranging from coalition building and role specialization to unexpected strategies, deception, and even cultural drift, all appearing without being explicitly programmed[6].

These developments expose the limits of traditional evaluation methods. Simple metrics like task completion or efficiency no longer capture the full spectrum of what multi-agent LLM systems can do or the risks they may pose.

Figure 1: Evolution of MAS Evaluation



3. CURRENT EVALUATION APPROACHES

3.1 Task-Specific Benchmarks

A number of evaluation environments have been developed for multi-agent systems, including ALFWorld[7], GAIA[8], and WebArena[9]. These settings are structured and relatively easy to replicate, which makes them useful for testing agents in controlled conditions. At the same time, their scope is narrow. They typically involve only a handful of agents, usually two to five, working toward simple, well-defined goals such as navigation or tool use.

3.2 Social Simulations

Other efforts take inspiration from social environments. Projects like Smallville[6] and various AI-driven town simulations model communities of agents interacting in open-ended ways. These simulations often produce fascinating emergent behaviors such as agents forming relationships, organizing gatherings, or even spreading gossip. However, the richness that makes them engaging also makes them difficult to measure systematically. Findings are often anecdotal, which limits their reproducibility and comparability.

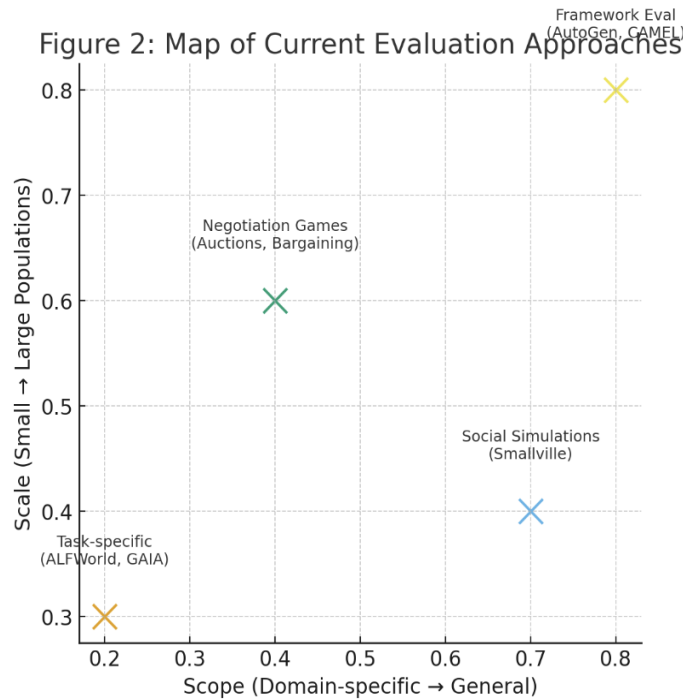
3.3 Negotiation and Cooperation Games

Structured games such as auctions, bargaining exercises, and resource-sharing scenarios have also been used to study collaboration between agents. These setups are well suited to exploring concepts like fairness, trade-offs, and economic efficiency. The drawback is that such games are highly stylized. They simplify real-world negotiation to the point where the results may not translate beyond the laboratory setting.

3.4 Framework-Specific Evaluations

Many agent frameworks, including AutoGen, CAMEL, and LangGraph, come with their own evaluation metrics. These usually track factors such as task success rates or the computational cost of running a

system. While these measures are valuable for improving a given framework, they remain fragmented. Because each framework uses its own standards, results are not easily comparable across systems.



4. TAXONOMY OF EVALUATION DIMENSIONS

To make sense of the diverse approaches in multi-agent evaluation, it is useful to group them into a set of core dimensions. These dimensions capture not only how well agents perform but also how resilient, interpretable, and scalable their collective behavior is. We outline five areas that together provide a more comprehensive lens for evaluation.

4.1 Collaboration Quality

One of the most straightforward measures is whether a group of agents can achieve its intended goal. Success rates provide a baseline, but deeper insight comes from looking at the efficiency of coordination, such as how much communication is required, how quickly tasks are completed, and whether resources like tokens or computing cycles are used effectively.

4.2 Emergent Behaviors

Multi-agent systems often display dynamics that were not explicitly designed. Agents may begin to adopt specialized roles, with one acting as a leader and another as a critic. They may experiment with different strategies or fall into recurring patterns of cooperation and conflict. Evaluating emergence means paying attention to whether these behaviors are stable and productive, or whether they result in cycles of confusion and breakdown.

4.3 Robustness and Safety

A reliable system should continue to function even when some agents behave unpredictably. This dimension considers resilience to malicious or selfish actors, as well as how errors spread through the network. It also includes the broader issue of trust and transparency: whether observers can understand the reasoning behind agent actions and whether outcomes remain consistent with intended objectives.

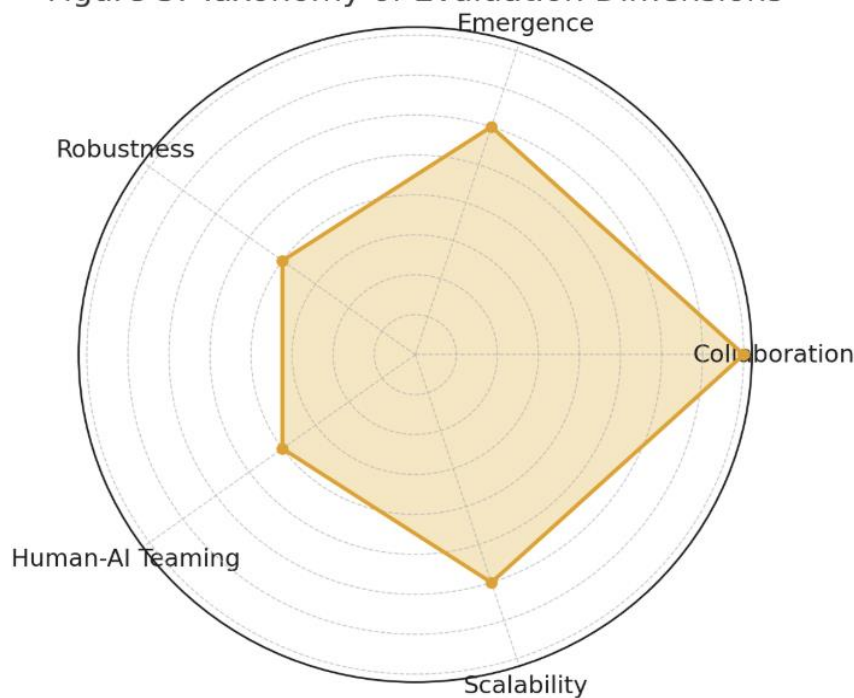
4.4 Human - AI Teaming

In many real applications, humans will not be replaced but instead embedded within multi-agent collectives. Evaluation here asks how effectively humans can oversee, guide, or collaborate with groups of agents. This involves examining whether communication remains clear, whether oversight becomes too demanding, and whether agent decisions stay aligned with human goals.

4.5 Scalability and Efficiency

Evaluation must also consider whether performance holds as the number of agents grows. A system that works smoothly with ten agents may struggle when expanded to one hundred or one thousand. Important factors include latency, computational cost, and environmental impact such as energy consumption. Scalability is both a technical and practical challenge, determining whether these systems can ever be applied in real-world settings.

Figure 3: Taxonomy of Evaluation Dimensions



5. GAPS IN CURRENT EVALUATION PRACTICES

5.1 Scale limitations

Most benchmarks are designed for small groups of agents, often no more than twenty. Real-world scenarios such as logistics or market simulations could easily involve hundreds or thousands, and current methods do not capture how systems behave at that scale.

5.2 Anecdotal emergence

Phenomena like role specialization, coalition-building, and even deception are frequently observed, but they are reported informally rather than measured systematically. Without reproducible methods, it is difficult to separate genuine emergence from isolated cases.

5.3 Fragmented benchmarks

In contrast to single-agent evaluation, which has common suites like GLUE or MMLU, multi-agent research lacks a unified benchmark. Each framework or study relies on its own standards, which makes comparison across systems nearly impossible.

5.4 Missing safety metrics

Risks such as collusion, manipulation, or the spread of misinformation are recognized, but there are few concrete tools or metrics to evaluate these issues in multi-agent environments.

5.5 Limited focus on human oversight

Most evaluations ignore the practical question of how humans can effectively supervise or guide groups of agents. Oversight at scale is likely to be critical in deployment, yet it remains largely untested.

5.6 Cross-domain disconnect

Research in fields like finance, robotics, and healthcare often proceeds independently, with little cross-pollination. This siloed approach prevents the development of evaluation practices that could generalize across domains.

6. FUTURE DIRECTIONS AND ROADMAP

6.1 Design principles for benchmarks

Future benchmarks should span multiple domains rather than focusing on narrow task environments. They need to scale to populations of hundreds of agents, since small groups cannot capture the dynamics of real-world applications. Equally important, benchmarks must be transparent and reproducible so that results can be trusted and compared across research groups.

6.2 Standardized metrics

To make evaluations consistent, the field needs shared measures. A Collaboration Index could capture how well agents balance success and efficiency against cost. A Trust Index could quantify alignment by weighing cooperative outcomes against deceptive or adversarial behaviors. An Emergence Score might measure the diversity of strategies alongside the degree of role specialization that develops within the system.

6.3 Large-scale synthetic societies

Building simulated environments such as digital cities, synthetic economies, or disaster response ecosystems would allow researchers to stress-test agent collectives in conditions closer to reality. These settings can reveal how governance structures form, how resources are distributed, and how resilient the system is under shocks.

6.4 Human-in-the-loop benchmarks

Since many deployments will involve humans working alongside or overseeing large groups of agents, benchmarks must reflect that interaction. Evaluation should include measures of interpretability, ease of oversight, and the burden placed on human supervisors when guiding swarms of agents.

6.5 Toward a Multi-Agent Benchmark Suite (MABS)

The long-term goal should be to create a gold-standard benchmark suite for multi-agent LLM systems, similar to how GLUE and MMLU transformed single-agent evaluation. Such a suite should be open-source, community-driven, and include a diverse set of tasks that reflect both technical challenges and societal considerations.

7. CONCLUSION

Multi-agent LLM systems mark an important step forward, moving beyond single-agent reasoning toward forms of collective intelligence. Despite this promise, evaluation practices are still fragmented and underdeveloped. In this review, we outlined the current landscape, identified key gaps, and proposed a taxonomy of dimensions that can guide more systematic benchmarking.

Looking ahead, the field would benefit from the development of a **Multi-Agent Benchmark Suite (MABS)** designed to capture collaboration quality, emergent behavior, robustness, human-AI interaction, and scalability. Standardized benchmarks of this kind would provide a shared foundation for progress. Without such efforts, research risks remaining anecdotal and siloed. With them, multi-agent systems have the potential to mature into safe, reliable, and widely deployable technologies.

REFERENCES:

1. Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D. and Steinhardt, J., 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
2. Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., Zhang, Y., Narayanan, D., Wu, Y., Kumar, A. and Newman, B., 2022. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*.

3. Srivastava, A., Rastogi, A., Rao, A., Shoeb, A.A., Abid, A., Fisch, A., Brown, A.R., Santoro, A., Gupta, A., Garriga-Alonso, A. and Kluska, A., 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on machine learning research*.
4. Samvelyan, M., Rashid, T., De Witt, C.S., Farquhar, G., Nardelli, N., Rudner, T.G., Hung, C.M., Torr, P.H., Foerster, J. and Whiteson, S., 2019. The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*.
5. Terry, J., Black, B., Grammel, N., Jayakumar, M., Hari, A., Sullivan, R., Santos, L.S., Dieffendahl, C., Horsch, C., Perez-Vicente, R. and Williams, N., 2021. Pettingzoo: Gym for multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 34, pp.15032-15043.
6. Park, J.S., O'Brien, J., Cai, C.J., Morris, M.R., Liang, P. and Bernstein, M.S., 2023, October. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology* (pp. 1-22).
7. Shridhar, M., Yuan, X., Côté, M.A., Bisk, Y., Trischler, A. and Hausknecht, M., 2020. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv preprint arXiv:2010.03768*.
8. Mialon, G., Fourrier, C., Wolf, T., LeCun, Y. and Scialom, T., 2023, October. Gaia: a benchmark for general ai assistants. In *The Twelfth International Conference on Learning Representations*.
9. Zhou, S., Xu, F.F., Zhu, H., Zhou, X., Lo, R., Sridhar, A., Cheng, X., Ou, T., Bisk, Y., Fried, D. and Alon, U., 2023. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*.