

A Reliability Control Framework for Robust Multi-Agent LLM Systems: Managing Workflows in Large Language Model Systems

Jasneet Arora^{1*} and Gurpreet Singh¹

^{1*}Department of Computer Science and Engineering, Chitkara University, Punjab, India.

Abstract

Large language model (LLM)-driven multi-agent systems are rapidly becoming integral to enterprise automation. Despite their increased use, real-world applications still face many reliability issues, including coordination problems, verification gaps, and recovery weaknesses. Recent empirical research calls for solutions to the problem that go beyond the management of tasks.

In this paper, we propose the AI Reliability Control Framework (ARCF), which incorporates seven layers, intended to enhance the reliability of production multi-agent large language model systems. This framework has a Reliability Control Layer that controls operations, identifies failure instances, and triggers automatic fault recovery mechanisms. ARCF has four main functionalities: (i) monitoring in real time according to the Multi-Agent System Failure Taxonomy (MAST), (ii) formally verified fault-recovery protocols with probabilistic reliability analysis, (iii) observability using OpenTelemetry GenAI standards, and (iv) safety verification through the Constitution for Artificial Intelligence framework.

Simulation evaluation of ARCF across various workflows used in enterprises reveals significant efficiency: it handled 89% of the failures. Additionally, real-time monitoring managed to encompass 93% of the current types of failures. These findings suggest that architectural supervision can help greatly improve large language model multi-agent systems in real-world use.

Keywords: Multi-Agent Systems, Large Language Models, AI Reliability, Fault Recovery, Observability, Constitutional AI

1 Introduction

Many special agents work together to perform processes that require multiple stages. Such systems are currently being applied extensively in many industries including fraud detection, automated customer service, finance analytics, and research support.

Multi-agent systems are more preferable to single-agent systems since they break down complex objectives into small tasks performed by specialized agents. It improves modularity, scalability, and adaptability to various types of workloads. Nevertheless, reliability remains the main problem with multi-agent systems application.

In practice, several challenges arise, including misusing tools, poor agent coordination, unexpected disruptions in the flow of workflows, and low-quality intermediate results. According to recent findings, the mentioned failures happen quite frequently in existing orchestration frameworks. Therefore, in addition to being related to modeling restrictions, reliability issues emerge because of architectural and implementation decisions. The above-mentioned insights highlight the necessity of developing a structured reliability engineering strategy that is tailored specifically to multi-agent LLM contexts.

Most of the existing orchestration systems tend to focus more on task orchestration and resource management. Although these aspects are significant, they usually depend on some kind of monitoring tools outside the system to fix failures that occur. Consequently, most orchestrations lack mechanisms to detect faults automatically and ensure workflow safety. Observability allows users to collect telemetry and trace data, while safety guardrails help filter outputs. On the other hand, such solutions work separately from one another without building a reliability solution together.

In order to fill this gap, this paper introduces the AI Reliability Control Framework (ARCF), an architecture aimed at enhancing the robustness of production multi-agent LLM systems. The framework has a specialized Reliability Control Layer that manages the execution stack and helps out in the event of failures. This layer brings together monitoring, failure recovery, observability, and safety checks in a coordinated design.

The main contributions of this paper are as follows:

- A seven-layer architecture that includes reliability supervision as a key component of multi-agent LLM systems.
- A monitoring engine that can detect known multi-agent failure modes in real time.
- Four fault-recovery algorithms that increase the probability of successful workflows through retries, fallback agents, checkpoint recovery, and circuit-breaker isolation.
- An observability and safety pipeline that follows emerging standards for LLM monitoring and responsible AI. A simulation-based evaluation shows major improvements in workflow completion and failure recovery.

The other sections of this paper is organised as follows. Section 2 reviews related work. Section 3 presents the proposed framework architecture. Section 4 describes the core algorithms and components. Section 5 outlines the experimental methodology. Section 6 reports results and analysis. Section 7 discusses limitations and future research directions. Section 8 concludes the paper.

2 Related Work

2.1 Multi-Agent LLM Orchestration Frameworks

The rapid development of large language models has led to the development of frameworks that support cooperation among multiple agents in order to complete complex tasks. The initial frameworks focused on allowing conversational agents to communicate with one another through structured conversations. These frameworks have further been expanded through subsequent developments in the areas of role specialization and workflow management.

Many of these orchestration frameworks are being often used in research and practice. AutoGen enables configurable conversational agents that collaborate through message passing and tool invocation. MetaGPT formalises standard operating procedures by assigning domain-specific roles to agents and structuring their interactions as assembly-line workflows. CAMEL introduced role-playing and cooperative planning mechanisms that encourage reflection and iterative improvement during task execution. LangChain remains a widely adopted infrastructure layer that provides integrations for tools, memory, and retrieval-augmented generation pipelines.

Comparative analyses reveal a notable disparity in reliability and task success among these frameworks. These differences show the importance of architectural design decisions, like coordination strategies, verification techniques, and error management. Though orchestration frameworks establish the foundation for multi-agent collaboration, they normally lack inherent features for reliability oversight.

2.2 Failure Taxonomy and Reliability in AI Agents

An analysis of failure reasons in multi-agent systems is a subject for intensive investigation. Recently published research results highlight recurring failure scenarios which occur in various modeling frameworks and application areas. There are several main failure types that can be classified as follows: problems with the design and definition of the multi-agent system, problems associated with inter-agent coordination, and problems related to task verification or termination.

A body of literature analyzing the dependability of AI agents proves that accuracy figures derived from a single trial cannot serve as a adequate performance indicator in practical settings. On the contrary, multi-trial assessment, fault injection, and stress testing are considered to be the appropriate ways to simulate actual working conditions. This demonstrates the necessity of having systems capable not only of controlling but also of monitoring tasks and detecting failures.

While the reliability issues become increasingly recognized, most of the current studies concentrate exclusively on the identification of failure types rather than proposing solutions to tackle them.

2.3 Observability and Monitoring for LLM Systems

With the adoption of LLMs into production environments, observability is becoming an indispensable need. With telemetry and tracing technologies, developers can

observe the system behaviour, find out any bugs that might occur and optimize performance. The OpenTelemetry tool is frequently used to help in collecting metric, log, and trace data. Also, new developments have created semantic standards made especially for LLM applications.

These conventions provide for standardized fields regarding model calls, tokens, latency, agents etc. As production monitoring and observability tools begin to support them, this allows the multiple agent interactions in production to be observed. However, these solutions are mainly used for data collection and visualization, and they cannot ensure reliability.

2.4 Safety and Guardrails for LLM Applications

Safety and alignment are among the key considerations when developing systems based on LLM. According to the Constitutional AI developed by Anthropic, one way that can be used in ensuring alignment is through the utilization of principles that would encourage a model to critique its output and make necessary changes. Other methods of achieving alignment include the application of guard rails and behavioral contract.

These are methods that ensure that the output is safe and free from errors but do not cover issues of safety and alignment at workflow levels, especially when there are multiple interacting agents.

2.5 Research Gap

The recent advances in research have led to substantial progress in orchestration models, fault analysis, observability, and safety guardrails. Nevertheless, the mentioned issues are rarely considered together. Current solutions do not provide an architecture with integrated monitoring, fault detection, automated recovery, observability, and safety validation.

Our paper fills this gap. We present an architecture that integrates all these components in a unified supervisory system for production-level multi-agent LLM architectures.

3 New Framework Architecture

The Artificial Intelligence Reliability Control Framework (ARCF) employs a layered approach that makes multi-agent language learning machines more reliable in production. This framework defines reliability management as an architectural element that includes monitoring, recovery, observability, and verification of safety as well.

3.1 Design Goals

This framework has been established with four primary objectives:

1. **End-to-end reliability:** The reliability features must operate throughout the entire workflow, encompassing task planning and the verification of the final output.
2. **Real-time failure awareness:** The system is required to regularly observe execution to find failures as they arise.

3. **Automated recovery:** If failures do occur, the system should be able to recover from these automatically without any manual intervention.
4. **Framework interoperability:** This architecture ought to function alongside existing orchestration frameworks without requiring a significant redesign of their foundational components.

These objectives help create of a supervisory layer which works independently of the execution stack.

3.2 Layered Architecture Overview

ARCF employs a seven-layer architecture. These layers include infrastructure as well as external interfaces and are designed to execute distinct functions in each layer.

1. **Foundation Model Layer:** This layer encloses the primary large language models and retrieval systems required by agents. It has both general-purpose foundation models and domain-specific knowledge sources. Model calls are arranged to capture metrics like latency, token usage, and error states.
2. **Data Operations Layer:** This layer manages both persistent and temporary system states. It has vector databases for retrieval-augmented generation, structured stores for intermediate workflow outputs, and event-based task ledgers that monitor the execution history. Concurrency controls help the state access to be always consistent, even when multiple agents work with shared resources.
3. **Agent Framework Layer:** This layer consists of agents responsible for planning, reasoning, and tool utilization. Agents work within reasoning-action loops and register their capabilities in an Agent Registry. Each agent has a dynamic health score, based on its long-term performance and reliability.
4. **Orchestration Layer:** The Orchestration Layer manages workflow execution. Tasks are organized as directed acyclic graphs (DAGs), where nodes define sub-tasks and edges describe dependencies. The orchestrator takes care of dependency resolution, agent selection, task dispatch, asynchronous communication among agents, and monitoring for timeouts and deadlocks. This layer is used as the execution backbone of the system.

3.3 Reliability Control Layer

The framework's main contribution is its Reliability Control Layer. It is a supervisory plane that observes system behavior, stepping in when failures occur. Unlike traditional orchestration components, this layer doesn't execute tasks directly. Instead, it monitors telemetry and coordinates recovery actions.

The layer has three primary functions:

1. **Failure detection:** Execution traces are used to recognize failure patterns, using a structured failure taxonomy.
2. **Recovery orchestration:** Detected failures are linked to suitable recovery strategies, like retries, selecting a fallback agent, rolling back checkpoints, or isolating an agent.

3. Reliability scoring: Agent performance is continuously evaluated. This improves the allotment of tasks based on agent's abilities and dependability.

This separation of reliability supervision from task execution means the framework can evolve its monitoring and recovery mechanisms independently.

3.4 Observability and Security Layer

The observability layer offers an organized perspective of the system's behavior. The execution events in the system are monitored using traces that are further categorized by session, workflow, agent, and calls to tools. This can be used in monitoring performance and failure analysis. In addition, there is runtime safety monitoring and prompt injection detection between agents.

3.5 Integration and API Layer

The uppermost layer within the architecture provides connectivity to external systems via its services interfaces. The services interface enables management of workflow configuration and monitoring through monitoring dashboard, reliability configuration, and integration with enterprise applications. As a result, this makes the architecture fit for deployment within production environments.

3.6 Formal Reliability Model

For measure the impact of recovery mechanisms, we provide workflow reliability using the probabilities of success of individual agents. Suppose there is a workflow consisting of n sequential agents. Then, $P(a_i)$ denotes the probability of successful completion of the task by agent a_i .

For analyzing the importance of recovery mechanisms, we introduce the concept of workflow reliability using the success probabilities of each individual agent. Assume that there is a workflow comprising n sequential agents. Then, $P(a_i)$ represents the probability of successful completion of the assigned task by agent a_i .

Without recovery mechanisms, workflow reliability is:

$$R_{\text{base}} = \prod_{i=1}^n P(a_i) \quad (1)$$

When we have the retry strategies, the effective success probability goes up, calculated as

$$P_{\text{retry}}(a_i) = 1 - (1 - P(a_i))^k \quad (2)$$

where k is the count of the retry attempts. When fallback agents are available, the composite probability further improves by incorporating their fallback success probability.

$$P_{\text{composite}}(a_i) = P_{\text{retry}}(a_i) + (1 - P_{\text{retry}}(a_i)) \cdot P(a_{\text{fallback}}) \quad (3)$$

This model demonstrates how recovery strategies can significantly boost end-to-end workflow reliability. It provides the analytical basis for the recovery algorithms described in the next section.

4 Major Components and Algorithms

This section studies the major components of the Reliability Control Layer and the algorithms that detect, mitigate, and recover from failures in multi-agent workflows.

4.1 Reliability-Aware Orchestrator

The orchestrator controls the flow of workflow execution through the use of a directed acyclic graph (DAG), which outlines all the tasks involved. Each node in this graph corresponds to a task, while edges denote the interdependencies between them. During the execution process, the orchestrator assigns agents according to their skills and scores for reliability. Afterward, the orchestrator executes the assigned tasks and keeps telemetry records. The system's control loop also identifies and rectifies any errors occurring during the entire process. Upon completion of a task, the orchestrator examines the corresponding telemetry data. If the results obtained are unacceptable, the task is forwarded to the recovery subsystem.

4.2 MAST-Aware Monitoring Engine

The monitoring engine is responsible for analyzing the execution telemetry and detecting failure patterns that emerge during the process. These include measuring factors such as response time, token consumption, format validation, tool output, and communication patterns.

The types of failures can be classified into three major categories:

1. **Specification/system design failures**, which encompass factors like timeout errors, context overflow, schema violations, and improper tool utilization.
2. **Misalignment between agents failures**, which involves factors like communication failures, output conflict, role violation, and workflow deadlock.
3. **Verification/termination failures**, which involve premature termination, infinite loop, and output quality degradation.

Failure events trigger the generation of an organized failure report that is forwarded to the recovery component.

4.3 Adaptive Exponential Backoff (Retry Strategy)

Transient failures often stem from temporary resource limitations, latency spikes, or intermittent tool errors. The framework addresses these by employing an adaptive exponential backoff retry strategy. With every attempt, the interval between retries increases exponentially. This interval subsequently adjusts according to a factor determined by the system's present success rate:

$$\text{AdaptFactor}(r) = \begin{cases} 2.0 & \text{if } r < 0.50 \\ 1.5 & \text{if } r < 0.75 \\ 1.0 & \text{otherwise} \end{cases} \quad (4)$$

4.4 Hierarchical Fallback Agent Selection

After all retries are used, the framework chooses another agent to handle the failed task. To find the most suitable one, it scores and ranks the available agents using the following function:

$$\text{score}(a_f) = w_1 \cdot H(a_f) + w_2 \cdot \text{CapScore}(a_f, t_i) + w_3 \cdot \text{HierBonus}(a_f) + w_4 \cdot \text{Cost}(a_f) \quad (5)$$

Here, $w_1 = 0.40$ (reliability), $w_2 = 0.30$ (capability), $w_3 = 0.20$ (hierarchy), and $w_4 = 0.10$ (cost). The execution context from the failed task then transfers to the fallback agent. This transfer prevents redundant recomputation and maintains workflow continuity.

4.5 Checkpoint-Based State Recovery

When failures affect multiple components or corrupt intermediate state, the framework relies on checkpoint-based recovery. The system saves snapshots of the workflow's state on regular basis. Should a severe failure occur, the framework checks if rolling back makes sense. It does this by comparing the estimated re-execution cost against the workflow's remaining value:

$$\text{rollback if: } \text{re_cost} \leq \rho \cdot \text{rem_val} \quad (6)$$

Here, ρ is a configurable cost-ratio threshold. If a rollback is justified, the system restores an earlier checkpoint and resumes execution without restarting the entire workflow.

4.6 Circuit Breaker Isolation

The agents which underperform consistently negatively impact the overall performance of the system. In order to reduce these repetitive failures, the framework applies a circuit breaker mechanism, that works in three states:

- **Closed:** normal operation.
- **Open:** a failing agent is temporarily isolated.
- **Half-open:** the recovery probe phase.

When an agent fails too often, the system temporarily pulls it from task allocation. After a recovery period, it runs lightweight probe tasks. If these confirm the agent's stability, the system then puts it back into service.

4.7 Verification and Safety Pipeline

Every output from the agents goes through a three-stage verification process:

- **Tier 1 – Schema validation:** This initial step confirms outputs match the expected structure.
- **Tier 2 – Safety assessment:** Safety checks are conducted here, drawing on both rule-based and principle-based methods, all in line with Constitutional AI principles.
- **Tier 3 – Logical consistency verification:** This final step confirms logical consistency, ensuring results make sense for the task and any prior outcomes.

If an output is not validated, it launches a retry, fallback, or rollback, on the basis of the particular failure. Together, these processes help automate identification and correction of problems, leading to a seamless recovery throughout.

5 Experimental Setup

This section details the evaluation of the proposed ARCF. The experiments were done to check reliability, performance, and safety under controlled fault conditions.

5.1 Evaluation Methodology

Assessing reliability in multi-agent systems necessitates more than just single-run accuracy metrics. To match real-world conditions, the system is tested multiple times with carefully introduced faults. This approach allows for the observation of failure detection and recovery actions across multiple trials.

A simulation of an enterprise business process was run using three agents working together to complete specific tasks. These tasks contained a Data Processing Agent that handled data intake and preprocessing, a Domain Analysis Agent that executed reasoning and analytical tasks, and a Report Generation Agent assigned with producing structured output. The agents work within a directed acyclic graph, dependent on one another.

To test robustness, failures were added to the system in different experiments at rates of 5%, 10%, and 20%. The 10% fault injection rate was chosen as the primary evaluation scenario, as it reflects realistic industrial conditions. The faults included random latencies and timeouts, execution errors of tools, context window overflow, agent crashes and delays, as well as faults in output generation.

Two configurations were assessed: (1) the standard orchestration process of the MAS without reliability algorithms; and (2) the fully operational ARCF system.

5.2 Evaluation Metrics

We checked the performance using many metrics, classifying them into four primary groups. In the performance category, we included Response Latency (P95), the 95th percentile of workflow completion time, and Observability Overhead, further delay introduced by telemetry instrumentation. For the reliability aspect, we kept track of the Workflow Completion Rate, Failure Recovery Rate, MAST Coverage Rate, Retry Success Rate, Fallback Success Rate, and Checkpoint Recovery Rate. We measured quality on the basis of Output Accuracy, which means the proportion of outputs that successfully met logical consistency criteria. Finally, our safety measure was the Safety

Block Rate, showing the percentage of outputs that were left out from release for safety reasons.

5.3 Implementation Details

We built the simulation environment around a modular orchestration prototype that generated synthetic workloads and telemetry traces. A monitoring engine then processed these traces in real time, triggering recovery algorithms whenever it detected failures. Our telemetry collection followed emerging semantic conventions for LLM observability. This method provided consistent measurement of latency, token use, and execution results. Each experimental setup was done several times to take care of the variations in agent behavior and the distribution of failures.

6 Results and Analysis

This section presents the experimental results and examines how the proposed framework affects workflow reliability, recovery effectiveness, and overall system performances.

6.1 Overall Reliability Improvements

Table 1 summarizes the comparative results between the baseline orchestration system and the ARCF-enabled system under a 10% fault injection rate.

Table 1 Comparative results between baseline and ARCF-enabled system under 10% fault injection rate

Metric	Baseline	ARCF
Workflow Completion Rate	78%	96%
Failure Recovery Rate	—	89%
MAST Coverage Rate	—	93%
Retry Success Rate	—	84%
Fallback Success Rate	—	79%
Checkpoint Recovery Rate	—	91%
Output Accuracy	—	—
Safety Block Rate	—	1.2%
Response Latency P95	1.8s	2.1s
Observability Overhead	—	3–4%

Introducing the Reliability Control Framework significantly improved workflow completion. The success rate rose to 96%, an 18 percentage-point gain over the baseline 78%. This shows the benefit of integrating embedding monitoring and automated recovery capabilities into the system design. The Reliability Control Framework managed to recover from failures in 89% of cases without human interaction.

6.2 Failure Detection Coverage

The suggested monitoring engine successfully identified 93% of recognized failure modes in real time. However, there were disparities in terms of detection among the various categories. The system achieved its best detection for specification and system design issues, and it also performed strongly with inter-agent misalignment failures. Task verification failures, though, had a slightly lower detection rate because subtle quality drops are tougher to spot. This shows that telemetry monitoring is able to find many of the structural problems, but assessing semantic quality remains challenging.

6.3 Recovery Algorithm Performance

Each recovery method helped in improving the overall reliability in a different way. The adaptive retry algorithm was able to successfully fix 84% of the temporary failures. This confirmed many multi-agent workflow issues are temporary, often solvable through controlled re-execution. Fallback activation worked in 79% cases, showing that critical tasks are able to take advantage from multiple capable agents. Checkpoint rollback achieved the highest success rate, restoring execution in 91% of severe failure scenarios. This highlighted the advantage of maintaining intermediate workflow state. The circuit breaker mechanism activated infrequently, isolating persistently failing agents in about 1.2% of executions. Despite its rare use, this mechanism kept repeated failures from spreading across workflows.

6.4 Performance Overhead

Reliability improvements were achieved with modest performance overhead. There was a small increase in Response latency experienced, now measuring 2.1 seconds (P95) compared to the previous 1.8 seconds. Additionally, the use of telemetry instrumentation contributed an extra 3–4% in latency. Most organizations see this a reasonable trade-off because reliability and accuracy are often more important than minor latency increases.

6.5 Safety Evaluation

Even with such measures in place, only 1.2% of the generated outputs were blocked owing to policy violations and unsafe material. Notably, even with such a high level of filtering being done, there was no significant decrease in the workflow rate.

6.6 Summary of Findings

Experiment results show that: integrated reliability supervision increases the probability of successful workflow execution; automated recovery procedures have the ability to fix most failure cases; observability and safety validation add only insignificant delay overhead; and recovery from checkpoints and retrying offer the biggest reliability boost. The results validate the effectiveness of our architectural proposal to ensure reliability for production multi-agent LLMs.

7 Discussion and Limitations

The results demonstrate that introducing a dedicated reliability supervision layer can substantially improve the robustness of multi-agent LLM workflows. The observed increase in workflow completion and the high failure recovery rate indicate that many production failures are not irrecoverable, but instead require timely detection and suitable mitigation strategies.

First, it should be noted that not all recovery techniques performed equally well in dealing with different types of problems. In particular, retry turned out to be very efficient when dealing with temporary issues, including increased latency and temporary problems with tools used for the job. The most beneficial effect could be achieved through the use of checkpointing to address complex problems impacting workflow status.

The other notable result was that the overhead associated with monitoring and telemetry was minimal latency-wise. The extra processing time required was fairly small, especially when taking into consideration the gains made in terms of reliability. This trade-off is very valuable in an enterprise setting, where accuracy and completion take priority over achieving minimal latency.

The evaluation also revealed a persistent challenge: detecting semantic quality degradation. While telemetry signals easily identify structural failures, automatically recognize minor reasoning flaws or content quality issues is much harder. Enhancing these semantic evaluation mechanisms ought to be a primary emphasis for forthcoming research.

Although the suggested framework is very reliable, it has a number of weaknesses as well. First of all, our study was conducted in a simulation environment that is somewhat similar to enterprise settings but might have additional complexity issues due to the different environment or external dependencies.

Furthermore, the monitoring component depends heavily on rule-based failure identification using telemetry information. This limits its effectiveness in detecting semantic and reasoning faults.

In addition, the checkpoint strategy incurs additional storage costs, which increase in proportion to the complexity of the workflow.

The future direction of our study will focus on incorporating learning-based failure classification, verification of recovery techniques and state transitions, generalization of the reliability model for parallel/asynchronous workflows, and the creation of adaptive agents that change behavior based on reliability experience.

8 Conclusion

This paper introduces the AI Reliability Control Framework (ARCF). It uses a multi-layer architecture to improve the reliability of multi-agent systems, especially those built with large language models. A key part of the framework is the Reliability Control Layer, which runs workflows, detects failures in real time, and starts automatic recovery when problems occur.

The proposed architecture brings together monitoring, fault recovery, observability, and safety verification in one unified design. Four specific recovery mechanisms,

adaptive retry, fallback agent selection, checkpoint-based rollback, and circuit breaker isolation, greatly improved workflow robustness when used in a reliability-aware orchestration pipeline.

Our experiments showed major improvements in operational reliability. Workflow completion hit 96%, a substantial increase over 78%. The system also achieved an 89% failure recovery rate, all while adding only modest latency. These findings indicate that prioritizing architectural reliability is essential for ensuring dependable production deployments of multi-agent LLM systems.

The next step is real-world application. Additionally, there will be a learning-based approach to detecting failures, and the framework will be expanded to handle more advanced processes. This framework can provide a concrete starting point for reliability engineering in future enterprise AI systems.

Acknowledgements. The author thanks the research community focused on multi-agent LLM reliability; their empirical foundations inspired this work.

Declarations

- **Funding:** Not applicable.
- **Conflict of interest:** The author has no competing interests to declare.
- **Ethics approval and consent to participate:** Not applicable.
- **Consent for publication:** Not applicable.
- **Data availability:** Not applicable (simulation-based study).
- **Code availability:** Not applicable.
- **Author contribution:** Jasneet Arora developed the framework, conducted the experiments, and wrote the manuscript.

References

- [1] Chase H (2022) LangChain. GitHub. <https://github.com/hwchase17/langchain>
- [2] Wu Q, Bansal G, Zhang J, et al (2023) AutoGen: Enabling next-gen LLM applications via multi-agent conversation. arXiv:2308.08155
- [3] Li G, Hammoud H, Itani H, et al (2023) CAMEL: Communicative agents for ‘Mind’ exploration of large language model society. In: Advances in Neural Information Processing Systems (NeurIPS), vol 36
- [4] Autono: A ReAct-based highly robust autonomous agent framework (2024)
- [5] Hong S, Zhuge M, Chen J, et al (2024) MetaGPT: Meta programming for a multi-agent collaborative framework. In: Proceedings of ICLR 2024. arXiv:2308.00352
- [6] Cemri M, Pan MZ, Yang S, et al (2024) Why multi-agent LLM systems fail. arXiv preprint

- [7] Drammeh P (2025) Multi-agent LLM coordination enables consistent high-quality decision assistance responding when something goes wrong. arXiv:2511.15755
- [8] Lewis P, Perez E, Piktus A, et al (2020) Retrieval-augmented generation for knowledge-intensive NLP tasks. In: Advances in Neural Information Processing Systems (NeurIPS), vol 33, pp 9459–9474
- [9] Tran HC, et al (2025) Toward engineering multi-agent LLMs: A protocol-driven approach
- [10] Chen M, et al (2025) Robustness of LLM-driven agent teams under member failures. In: ICLR Workshop. OpenReview:bkiM54QftZ
- [11] Anonymous (2026) Towards a science of AI agent reliability. arXiv:2602.16666
- [12] Yao S, Zhao J, Yu D, et al (2023) ReAct: Synergizing reasoning and acting in language models. In: Proceedings of ICLR 2023. arXiv:2210.03629
- [13] Wei J, Wang X, Schuurmans D, et al (2022) Chain-of-thought prompting elicits reasoning in large language models. In: Advances in Neural Information Processing Systems (NeurIPS)
- [14] OpenTelemetry (2024) LLM observability with OpenTelemetry. <https://opentelemetry.io/blog/2024/llm-observability/>
- [15] OpenTelemetry (2025) AI agent observability. <https://opentelemetry.io/blog/2025/ai-agent-observability/>
- [16] Datadog (2025) LLM observability and OTel GenAI semantic conventions. <https://www.datadoghq.com/blog/llm-otel-semantic-convention/>
- [17] ALAS: Adaptive LLM agent scheduling with fault-tolerant compensators (2025)
- [18] Fault-tolerant sandboxing for AI coding agents using transactions (2025)
- [19] Anonymous (2025) Open challenges in multi-agent security. arXiv:2505.02077
- [20] Li J, et al (2020) Chaos engineering for distributed systems. IEEE Software 37(5):36–43
- [21] Bai Y, Jones A, Ndousse K, et al (2022) Constitutional AI: Harmlessness from AI feedback. arXiv:2212.08073
- [22] Rebedea T, Dinu R, Sreedhar M, et al (2023) NeMo guardrails: A toolkit for controllable and safe LLM applications. In: Proceedings of EMNLP System Demonstrations

- [23] Anonymous (2026) Agent behavioral contracts: Formal specification and runtime enforcement. arXiv:2602.22302
- [24] Dong Y, et al (2025) A survey on large language model safety. Artificial Intelligence Review. <https://doi.org/10.1007/s10462-025-11389-2>
- [25] Park JS, O'Brien J, Cai CJ, et al (2023) Generative agents: Interactive simulacra of human behavior. In: Proceedings of UIST 2023. ACM
- [26] OpenAI (2023) GPT-4 technical report. arXiv:2303.08774
- [27] Bommasani R, Hudson DA, Aditi E, et al (2021) On the opportunities and risks of foundation models. arXiv:2108.07258
- [28] National Institute of Standards and Technology (2023) Artificial Intelligence Risk Management Framework (AI RMF 1.0). NIST