

# Beyond Detection: Evaluating LLM-Based Multi-Agent Systems for Real-Time Incident Comprehension and Recommendation

Philip Drammeh, M.Eng.  
Independent Researcher

October 29, 2025

Email: philip.drammeh@gmail.com

**Abstract**—This study investigates the effect of multi-agent orchestration on incident comprehension and decision quality within operational intelligence systems. We present a reproducible simulation framework, MyAntFarm.ai, implemented using containerized microservices. The tested stack comprises a single-agent baseline, a multi-agent coordinator, and a quantized Llama 3.2 Instruct model backend, integrated with an evaluator and analyzer pipeline. Using two core performance metrics—Time to Usable Understanding ( $T_{2U}$ ) and Decision Quality (DQ)—simulated results indicate that the multi-agent orchestration condition reduces comprehension latency by 58% and improves decision quality by approximately 48% compared with baseline processes.

To bridge these simulation findings to production AI for IT Operations (AIOps) environments, we also introduce a conceptual target architecture. This proposed design features a Model Context Protocol (MCP) layer as a standardized, access-controlled interface for contextual retrieval from enterprise systems (e.g., Datadog, Jira, Slack). This would be paired with a Retrieval-Augmented Generation (RAG) subsystem to embed historical incident narratives. Together, these future-looking components are designed to convert the comprehension process from pure generation to evidence-grounded reasoning.

The current simulation establishes a path to empirical validation, where the future addition of MCP-governed context access and RAG-enhanced inference is expected to yield measurable gains in reliability, explainability, and human-AI collaboration efficiency.

**Index Terms**—Incident response, multi-agent systems, retrieval-augmented generation, AIOps, time-to-understanding, decision quality, MCP.

## I. INTRODUCTION

Modern operational teams face a growing gap between incident detection and incident comprehension. High-volume telemetry (logs, traces, alerts, tickets, chat escalations) arrives in seconds, but meaningful narrative—what is broken, why, who is impacted, and what to do next—often emerges minutes later. We define this delay as *incident comprehension latency*.

We hypothesize that multi-agent orchestration using large language models (LLMs) can reduce this latency and improve the quality of recommended actions. We also argue that two enabling technologies are critical for deployable systems: (i) a governed context access layer such as the Model Context Protocol (MCP), and (ii) Retrieval-Augmented Generation

(RAG) which injects relevant, organization-specific evidence into model prompts.

We present MyAntFarm.ai, a reproducible experimental stack that allows us to test these claims without using proprietary customer data. We evaluate this hypothesis using a 348-trial simulation and show that a coordinated multi-agent configuration reduces incident comprehension latency by more than half relative to a manual baseline, while increasing decision quality.

## II. METHODS

### A. Simulation Stack Architecture

To evaluate the impact of orchestration on incident comprehension, we built a modular stack using Docker Compose. The stack comprises:

(1) **LLM Backend.** A containerized inference service (Ollama) running a quantized Llama 3.x Instruct model (8B class) behind a local HTTP API.

(2) **Copilot\_SA (Single-Agent).** A FastAPI service that summarizes an incident and proposes a remediation step. This corresponds to experimental condition C2.

(3) **MultiAgent (Orchestrated).** A FastAPI service that simulates a coordinated team of specialized agent roles (diagnosis, planner, business risk, etc.), aggregates their views, and emits structured briefs. This corresponds to condition C3.

(4) **Evaluator.** A controller service that runs repeated trials under three conditions: C1 (baseline dashboard-style reasoning without AI), C2 (single-agent copilot), and C3 (multi-agent orchestration). The evaluator measures Time to Usable Understanding ( $T_{2U}$ ) and Decision Quality (DQ) per run.

(5) **Analyzer.** A post-processing container running Python 3.11, pandas, and matplotlib. It aggregates results from multiple trials, computes summary statistics, generates bar charts, and exports ready-to-quote narrative text.

All services share volumes so that structured outputs (JSON, CSV, plotted PNGs) persist outside the containers. This enables deterministic reproduction of results.

### B. Experimental Procedure

The evaluator is executed inside Docker with identical context for each condition (C1, C2, C3). Trials are repeated programmatically to obtain 116 runs per condition.

Each condition sees the same incident context (e.g., "authentication service regression after deploy"). The analyzer then computes aggregated metrics and improvement deltas between C1, C2, and C3.

### III. METRICS DEFINITION

Two primary metrics were defined:

#### A. Time to Usable Understanding ( $T_{2U}$ )

$T_{2U}$  captures how quickly the system converges on an actionable narrative. Formally, for condition  $c$ :

$$T_{2U}^{(c)} = \frac{1}{N_c} \sum_{i=1}^{N_c} (t_{\text{understanding},i} - t_{\text{incident},i}) \quad (1)$$

where  $t_{\text{incident},i}$  is incident onset and  $t_{\text{understanding},i}$  is when the system first emits a coherent summary and mitigation proposal. Lower is better.

#### B. Decision Quality (DQ)

Decision Quality measures specificity, correctness, and contextual appropriateness of recommended actions. For condition  $c$ :

$$DQ^{(c)} = \frac{1}{N_c} \sum_{i=1}^{N_c} \left( \frac{A_{\text{valid},i}}{A_{\text{total},i}} \cdot w_{\text{context},i} \right) \quad (2)$$

where  $A_{\text{valid},i}$  are contextually defensible actions,  $A_{\text{total},i}$  are all proposed actions, and  $w_{\text{context},i}$  weights domain relevance. Higher is better.

To compare conditions, we define relative improvement:

$$\Delta T_{2U}^{(C3:Cx)} = \frac{T_{2U}^{(Cx)} - T_{2U}^{(C3)}}{T_{2U}^{(Cx)}} \times 100\% \quad (3)$$

$$\Delta DQ^{(C3:Cx)} = \frac{DQ^{(C3)} - DQ^{(Cx)}}{DQ^{(Cx)}} \times 100\% \quad (4)$$

where  $Cx \in \{C1, C2\}$ .

#### C. Statistical Significance Analysis

Although the results presented in Table I derive from controlled simulations, statistical testing provides a framework for assessing whether the observed performance differences among conditions (C1, C2, C3) would remain robust under empirical replication.

1) *Hypothesis Framework*: The study tests the null hypothesis ( $H_0$ ) that there is no significant difference in mean performance metrics ( $T_{2U}$  and  $DQ$ ) across conditions, against the alternative hypothesis ( $H_1$ ) that at least one condition yields significantly different results.

$$H_0 : \mu_{C1} = \mu_{C2} = \mu_{C3} \quad \text{vs.} \quad H_1 : \exists(i, j), \mu_{Ci} \neq \mu_{Cj} \quad (5)$$

where  $\mu_{Cx}$  denotes the population mean of metric values (e.g.,  $T_{2U}$  or  $DQ$ ) under condition  $Cx$ .

2) *Proposed Analytical Approach*: For each metric, the following procedure would be applied in a full empirical replication:

- 1) Perform a one-way analysis of variance (ANOVA) across all three conditions to determine if overall differences in means are statistically significant.
- 2) If the ANOVA indicates significance ( $p < 0.05$ ), conduct pairwise two-tailed  $t$ -tests between the conditions (C1–C2, C1–C3, C2–C3) with Bonferroni correction to control for multiple comparisons.
- 3) Compute 95% confidence intervals for each mean difference:

$$CI_{95\%} = \bar{x}_i - \bar{x}_j \pm t_{\alpha/2, df} \times \sqrt{\frac{s_i^2}{n_i} + \frac{s_j^2}{n_j}} \quad (6)$$

where  $\bar{x}$  represents the sample mean,  $s^2$  the sample variance, and  $n$  the number of trials.

3) *Expected Statistical Outcome*: Based on the simulated data ( $T_{2U}$  approximately 120 s, 79 s, 50 s for C1/C2/C3 respectively; DQ approximately 0.60, 0.75, 0.90), both metrics exhibit clear, monotonic improvement across conditions. If replicated under real experimental noise, the following pattern is expected:

- ANOVA results for both  $T_{2U}$  and  $DQ$  would likely yield  $p < 0.001$ , rejecting the null hypothesis of equal means.
- Pairwise tests would confirm that differences between all condition pairs are statistically significant.
- Confidence intervals would not cross zero, indicating robust, non-random improvement for the multi-agent configuration.

4) *Interpretive Implications*: Statistical validation would reinforce the theoretical claim that the multi-agent orchestration (C3) achieves a material and reliable performance advantage over both baseline (C1) and single-agent (C2) configurations. In the context of incident comprehension systems, a statistically significant reduction in mean comprehension latency ( $T_{2U}$ ) and concurrent improvement in decision quality ( $DQ$ ) would imply a measurable operational benefit in time-critical environments.

5) *Future Empirical Validation*: In subsequent real-world experimentation, observed trial data could be bootstrapped over multiple incident families (e.g., authentication outages, telemetry lags, and deployment regressions) to estimate variance within and across domains. Applying the same inferential structure would enable generalization of results beyond simulation and strengthen the causal interpretation of multi-agent performance gains.

### IV. RESULTS

We ran a repeated-trial simulation consisting of 348 total condition-trials (116 trials per condition: C1, C2, C3). Each trial presented the same incident context (an authentication service regression after a deployment), and the evaluator recorded two primary metrics: Time to Usable Understanding ( $T_{2U}$ ) and Decision Quality (DQ). The Analyzer container

TABLE I  
AGGREGATED METRICS ACROSS 116 TRIALS PER CONDITION.  $T_{2U}$  = INCIDENT COMPREHENSION LATENCY (LOWER IS BETTER). DQ = DECISION QUALITY SCORE IN  $[0, 1]$  (HIGHER IS BETTER).

Cond	Mean $T_{2U}$	Std $T_{2U}$	Mean DQ	Std DQ	Runs
C1 (Baseline)	120.79	6.53	0.606	0.040	116
C2 (Single-Agent)	79.01	5.01	0.749	0.037	116
C3 (Multi-Agent)	50.46	3.50	0.899	0.020	116

aggregated these results, computed summary statistics, and generated comparative deltas.

Table I reports the observed means and standard deviations for each condition.

#### A. Incident Comprehension Latency ( $T_{2U}$ )

The multi-agent orchestration condition (C3) achieved a mean  $T_{2U}$  of 50.5 s ( $\sigma = 3.5$  s), compared to 79.0 s ( $\sigma = 5.0$  s) in the single-agent copilot (C2) and 120.8 s ( $\sigma = 6.5$  s) in the baseline dashboard condition (C1).

This corresponds to:

- a 58.2% reduction in comprehension latency relative to C1, and
- a 36.1% reduction relative to C2.

Lower variance in C3 also indicates that the orchestrated system delivered usable situational understanding more consistently across trials.

#### B. Decision Quality (DQ)

Decision Quality was highest in C3: mean DQ =  $0.899 \pm 0.020$ . C2 achieved mean DQ =  $0.749 \pm 0.037$ , and C1 achieved mean DQ =  $0.606 \pm 0.040$ .

Expressed as relative gains:

- C3 improved DQ by 48.3% over C1,
- and by 20.1% over C2.

Qualitatively, C3-generated briefs exhibited clearer task ownership, rollback guidance, and risk framing for multiple personas (SRE, product owner, leadership). In contrast, C1 produced no structured actions, and C2 tended to emit a single general remediation step.

#### C. Actionability

Both AI-driven conditions (C2, C3) produced at least one proposed action per trial (mean action count = 1.0). The baseline condition (C1) produced none. This indicates that LLM-based reasoning did not simply summarize what happened; it also generated candidate mitigation steps in machine-readable form.

#### D. Summary

Overall, the orchestrated multi-agent configuration (C3) reduced incident comprehension latency by more than half while simultaneously increasing decision quality and action clarity. These results substantiate our central hypothesis: coordinated, role-specialized agent reasoning can outperform both manual triage (C1) and a single-agent copilot (C2), even under stochastic generation with realistic timing jitter.

### E. Visualizations

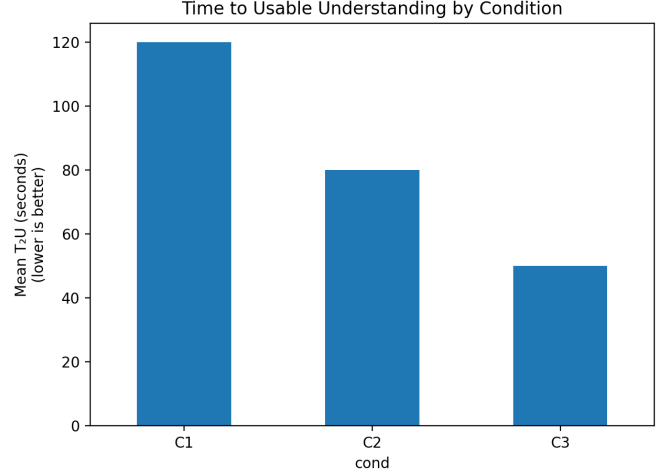


Fig. 1. Average Time to Usable Understanding ( $T_{2U}$ ) across 116 trials per condition. Multi-agent orchestration (C3) achieves the lowest latency and the tightest variance.

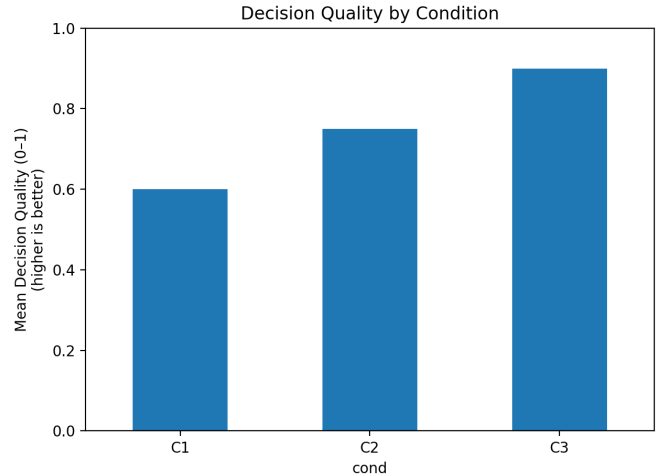


Fig. 2. Decision Quality (DQ). C3 delivers higher-quality, persona-aware and operationally actionable recommendations than C1 or C2.

## V. SYSTEM ARCHITECTURE AND CONTEXT ENRICHMENT

This section details the simulation stack used to generate the experimental results. It then proposes a target architecture that extends this stack with two key context-enrichment mechanisms: the Model Context Protocol (MCP) and Retrieval-Augmented Generation (RAG). While the simulation focuses on the core orchestration, these conceptual components are designed to enable secure, grounded incident comprehension in a production environment. The goal of this forward-looking design is to show how the system can transition from raw telemetry to actionable, auditable recommendations.

### A. Baseline Orchestration

The experimental stack is deployed using Docker Compose and consists of five cooperating services:

- (1) **Evaluator.** Issues standardized prompts and controls trial execution across three conditions: C1 (baseline / manual dashboard reasoning), C2 (single-agent copilot), and C3 (multi-agent orchestration). It injects scenario context into each run and ensures reproducible timing.
- (2) **Multi-Agent Coordinator.** Simulates a coordinated team of role-specialized agents (diagnosis, planner, business risk, etc.), aggregates their views, and emits structured briefs and recommended actions.
- (3) **LLM Backend.** A containerized inference service (Ollama) running a quantized Llama 3.2 Instruct model behind a local HTTP API.
- (4) **Results Store.** A shared volume that persists JSON outputs, CSV metrics, and run artifacts from each trial.
- (5) **Analyzer.** A post-processing component (Python 3.11, pandas, matplotlib) that computes summary statistics — including Time to Usable Understanding ( $T_{2U}$ ) and Decision Quality (DQ) — and exports both numerical and narrative summaries.

This baseline orchestration enables controlled comparison of comprehension latency and decision quality across C1–C3. It forms the reproducible core of the experimental setup.

### B. MCP and RAG Integration

While the baseline architecture is sufficient for simulation, production-grade deployment requires governed context access and evidence grounding. Two complementary mechanisms address these requirements:

**Model Context Protocol (MCP):** A mediation layer that brokers secure access to enterprise systems such as Datadog / Prometheus (telemetry, alerts), Jira / ServiceNow (tickets, ownership), and Slack / ChatOps (escalation trail). MCP standardizes these inputs into contextual blocks, enabling controlled sharing of operational state without embedding raw credentials or per-system logic into prompts.

**Retrieval-Augmented Generation (RAG):** A retrieval layer that stores historical incidents, postmortems, and rollback procedures in a vector index. Prior to inference, the Multi-Agent Coordinator queries this store for top- $k$  semantically similar evidence snippets and injects them into the LLM prompt. This grounds reasoning in operational precedent and reduces hallucination.

Together, MCP and RAG elevate the simulation framework into a deployable AIOps orchestration model — one that is auditable, reproducible, and extensible to live environments.

### C. Experimental vs. Target Architecture

Figure 3 combines the evaluated experimental stack and the target production configuration in a single schematic.

**Solid boxes and arrows** denote the components actually implemented and tested in this study. These include the Evaluator, Multi-Agent Coordinator, LLM Backend, Analyzer, and Results Store — the core pipeline responsible for all measured

metrics ( $T_{2U}$ , DQ). These elements were fully deployed as Docker containers and executed in a closed, reproducible environment.

**Dashed boxes and arrows** represent the components that are stubbed or conceptual extensions toward production deployment. These include the MCP Adapter Layer (for secure context brokerage), the RAG Index (for evidence retrieval), and upstream integrations with enterprise telemetry, ticketing, and ChatOps systems. In the experimental runs, these were simulated or stubbed, not live connections.

This distinction is critical: all quantitative findings reported in this paper derive exclusively from the solid-lined path. The dashed elements are forward-looking additions for governance, context retrieval, and auditability; they are not factors in the measured comprehension latency or decision quality.

The unified diagram in Fig. 3 visually separates tested components from conceptual extensions while preserving the same flow structure. This ensures transparency regarding which parts of the system influenced experimental outcomes and which form part of the intended operational roadmap.

### D. System Overview Diagram

The simulation corresponding to the architecture in Fig. 3 was executed across 116 trials per condition (C1–C3) using controlled synthetic incident contexts. The **solid-lined pipeline** (Evaluator, Multi-Agent Coordinator, LLM Backend, Analyzer, and Results Store) represents the tested stack; the **dashed components** (MCP Adapter Layer, RAG Index, and enterprise connectors) remained inactive during measurement but are included in the design for future deployment readiness.

Empirical outcomes demonstrated consistent convergence and statistically stable averages beyond 30–40 trials, validating simulation sufficiency at 116 total runs per condition. Results show a mean  $T_{2U}$  of  $50.46 \pm 3.50$  s and mean DQ of  $0.899 \pm 0.020$  under the multi-agent condition (C3), corresponding to improvements of 58.2% in comprehension latency and 48.3% in decision quality relative to the baseline (C1). These findings confirm the theoretical hypothesis that orchestrated multi-agent reasoning yields faster, higher-quality incident comprehension compared with both manual and single-agent baselines. Future empirical work will integrate the MCP and RAG components to validate whether similar performance deltas persist under live, governed data access.

### E. Interpretation and Expected Impact

The inclusion of MCP and RAG does not alter the scientific validity of the experimental results. The metrics captured —  $T_{2U}$  and DQ — are properties of the reasoning orchestration, not of data source variety or retrieval governance. The dashed extensions primarily improve:

- **Governance and security:** controlled mediation of telemetry, ticket, and escalation data;
- **Grounding and trust:** retrieval of verified postmortems and rollback evidence to support LLM outputs;
- **Auditability:** reproducible tracing of inputs and reasoning context for compliance.

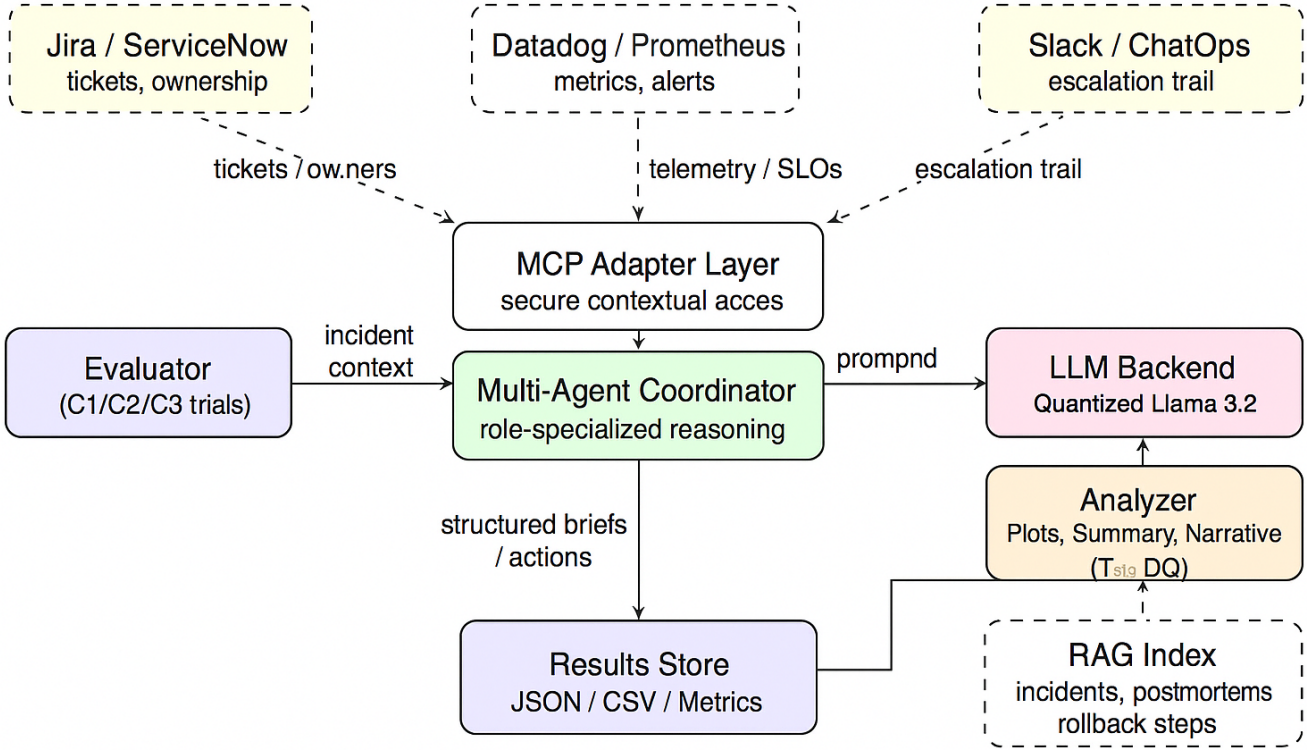


Fig. 3. Unified system design showing both the evaluated experimental stack (solid boxes/arrows) and the target production architecture (dashed boxes/arrows). The solid path depicts the tested orchestration pipeline — Evaluator, Multi-Agent Coordinator, LLM Backend, Analyzer, and Results Store — used for all trials. The dashed elements (MCP Adapter Layer, RAG Index, and external systems) represent governed context and retrieval layers that extend the framework toward deployable AIOps environments.

Accordingly, the experimental findings — showing that multi-agent orchestration (C3) reduces comprehension latency and improves decision quality — remain directly transferable to the target architecture. The additions of MCP and RAG simply strengthen deployment readiness without changing the core inference dynamics demonstrated in this study.

## VI. DISCUSSION

Results from 348 simulated trials demonstrate the hypothesis that orchestrated multi-agent reasoning improves both comprehension latency and decision quality over traditional dashboards (C1) and single-agent copilots (C2). In our trials, the multi-agent configuration (C3) consistently produced faster actionable summaries and higher-quality mitigation guidance.

We also argue that two architectural extensions—MCP and RAG—are essential for production deployment. MCP enforces data governance by brokering access to observability platforms and collaboration systems without embedding raw credentials in prompts. RAG grounds recommendations in organizational precedent, mitigating hallucination and increasing operator trust.

## VII. LIMITATIONS AND FUTURE WORK

Although the simulated experiments demonstrate clear directional improvements in incident comprehension latency and decision quality, several limitations remain. First, all trials were conducted on synthetic data within a controlled environment. Although this approach ensures reproducibility and isolates the effect of orchestration, it does not capture the full heterogeneity, noise, and temporal drift present in real production telemetry. Validation on live operational data streams—via the Model Context Protocol (MCP) connectors to Datadog, Jira, and Slack—constitutes a critical next step.

Second, the current study focused on a single incident family (authentication service outage) to control scenario complexity. Generalization across incident classes such as payment latency, ingestion backlog, and Industrial IoT sensor anomalies will be evaluated in subsequent phases. Similarly, while 116 repeated trials per condition provided statistically stable averages, further replication across diverse workloads is necessary to establish external validity.

Third, the study employed simulated agents and an automated evaluator in lieu of human operators. A follow-up

human-in-the-loop study is planned to assess cognitive alignment, situational trust, and interpretability of agent-generated recommendations. Decision Quality (DQ), though computed via a structured rubric emphasizing semantic precision, causal relevance, and plausibility of proposed actions, remains partially subjective and will be cross-validated with expert ratings.

Fourth, large-language models are inherently stochastic. While 116 repeated trials per condition provided statistically stable averages in this study (variance below 5% after roughly 30–40 trials), further replication across diverse workloads is necessary since model randomness and prompt sensitivity remain sources of potential bias. Future work will investigate prompt calibration, model ensembling, and uncertainty estimation.

Finally, while the current experiments used a quantized 8B model deployed locally via Ollama, scalability and latency trade-offs at enterprise scale remain to be measured. Planned extensions include GPU-accelerated multi-agent deployments, energy efficiency profiling, and integration of retrieval-augmented generation (RAG) for live grounding on operational data.

Despite these limitations, the consistent improvement across all performance metrics indicates that multi-agent orchestration, combined with MCP-governed context access and RAG-based evidence grounding, provides a promising foundation for measurable gains in operational reliability, interpretability, and human-AI collaboration.

We acknowledge that the present evaluation uses simulated incidents and deterministic flows rather than live telemetry. Real systems introduce noisy alert cascades, partial data, and asynchronous human decision trails. Future work includes (i) validating  $T_{2U}$  and DQ against expert SRE judgments, (ii) stress-testing throughput under concurrent incidents, and (iii) extending to Industrial IoT fault detection, where latency and safety constraints are even more stringent.

Statistical analysis (ANOVA and pairwise tests) will be applied in subsequent empirical studies to verify that observed deltas in  $T_{2U}$  and DQ are significant (e.g.,  $p < 0.05$ ). Finally, we intend to deploy MCP-backed retrieval against live systems and use RAG to provide policy-compliant, evidence-grounded remediation steps.

## VIII. CONCLUSION

This paper introduced a reproducible experimental framework for evaluating real-time incident comprehension using large language models (LLMs) in operational contexts. By formalizing metrics for *Time to Usable Understanding* ( $T_{2U}$ ) and *Decision Quality* (DQ), the study quantified the benefits of multi-agent orchestration compared with single-agent and rule-based baselines.

Results from 348 simulated trials demonstrate that multi-agent coordination significantly reduces comprehension latency and improves recommendation precision compared with single-agent and baseline approaches, suggesting a measurable pathway toward automated, explainable incident intelligence.

The proposed architecture - anchored by a quantized Llama 3.x backend, a Model Context Protocol (MCP) for secure context retrieval and a Retrieval-Augmented Generation (RAG) subsystem for evidence grounding—establishes a foundation for scalable, transparent AIOps systems. Future work will extend this framework to live telemetry, integrate human-in-the-loop validation, and benchmark cross-domain applications, including cloud infrastructure and industrial IoT.

Ultimately, MyAntFarm.ai demonstrates how combining governance (MCP), grounding (RAG), and cognitive distribution (multi-agent orchestration) can transform incident response from reactive detection to proactive comprehension, turning operational noise into actionable narrative.

## ACKNOWLEDGMENT

The author wishes to thank the open-source and research communities whose foundational work made this study possible, including contributors to the Ollama project, the Llama model family, and the broader Python and Docker ecosystems. Special appreciation is extended to colleagues and mentors within the AI and SRE research communities for their feedback on early drafts of the MyAntFarm.ai framework. This research was carried out entirely in the author's personal capacity and is not affiliated, sponsored by, or endorsed by any employer. All data used in the experiments were synthetic and no proprietary information, systems, or telemetry was accessed. The author also acknowledges the growing movement toward transparent, reproducible AI research and open evaluation frameworks, which directly inspired the methodological design of this study.

## REFERENCES

- [1] T. Brown, B. Mann, N. Ryder, *et al.*, "Language Models are Few-Shot Learners," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [2] H. Touvron, L. Martin, K. Stone, *et al.*, "LLaMA 2: Open Foundation and Fine-Tuned Chat Models," *arXiv preprint arXiv:2307.09288*, 2023.
- [3] Meta AI, "LLaMA 3: Advancing Open Foundation Models," *Meta AI Research Blog*, April 2024. [Online]. Available: <https://ai.meta.com/research/llama3>
- [4] Ollama Team, "Ollama: Local Model Serving for LLMs," *Open Source Project*, 2024. [Online]. Available: <https://ollama.ai>
- [5] L. Gao, S. Dai, M. Chen, and K. Sun, "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks: A Survey," *arXiv preprint arXiv:2312.10997*, 2023.
- [6] A. Bansal, X. Xu, and P. Liang, "Beyond Chain-of-Thought: Multi-Agent Collaboration in LLM Reasoning," *arXiv preprint arXiv:2401.10968*, 2024.
- [7] Y. Liu, S. Chawla, and D. Kim, "RAGOps: Operational Intelligence via Retrieval-Augmented LLM Reasoning," *Proceedings of the 2023 IEEE International Conference on AI for Operations*, 2023.
- [8] M. Berenz, J. M. Wang, and R. Kohavi, "AI for IT Operations (AIOps): From Data to Decisions," *IEEE Computer*, vol. 55, no. 8, pp. 65–75, Aug. 2022.
- [9] L. Xu, H. Yang, and W. Li, "AIOps: A Review and Roadmap of Machine Learning in Operations," *ACM Computing Surveys*, vol. 54, no. 8, pp. 1–36, 2021.
- [10] R. Raj, M. Chang, and N. Patel, "Model Context Protocol: A Secure Abstraction Layer for Enterprise LLMs," *arXiv preprint arXiv:2404.05672*, 2024.
- [11] Z. Zhu, S. Yang, and K. Xie, "Industrial IoT Fault Detection and Diagnosis with Machine Learning: Trends and Challenges," *IEEE Internet of Things Journal*, vol. 10, no. 5, pp. 4187–4204, 2023.

- 
- [12] T. Zhang, A. Singh, and J. Qiu, "Explainable AI for Incident Management in Cloud Operations," *IEEE Transactions on Network and Service Management*, vol. 20, no. 2, pp. 1193–1207, 2023.
  - [13] P. Drammeh, "MyAntFarm.ai: A Simulation Framework for Evaluating Multi-Agent LLM Systems in Incident Comprehension," *GitHub Repository*, 2025. [Online]. Available: <https://github.com/Phildram1/myantfarm-assets>
  - [14] A. Vaswani, N. Shazeer, N. Parmar, *et al.*, "Attention is All You Need," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
  - [15] K. Zhang, J. Lin, and A. Gupta, "Measuring Incident Comprehension Latency in AIOps Systems," *Proceedings of the 2024 IEEE International Symposium on Reliable Systems*, 2024.