# Architectures For Safety And Scalability In Agentic AI Systems

**Dr. Dipak Kadve[1] 0000-0002-2571-9973   Dr. Binod Kumar[2]   0000-0002-6172-7938**
**Prof. Vishal Gejge[3] 0009-0005-7870-3974**
**JSPM's Rajarshi Shahu College of Engineering, MCA Dept. Pune Maharashtra India**

**Abstract**

Agentic AI—systems of autonomous agents that perceive, reason, plan and act—are rapidly moving from proof-of-concept to production use in enterprise and research settings. This paper proposes a principled framework for **safe and scalable multi-agent systems** (MAS) that combines (1) specialized small agents for high-throughput repetitive tasks, (2) a coordinator/planner agent for long-horizon goal decomposition, and (3) an explicit **Human-in-the-Loop (HITL)** oversight and intervention layer. We present a hierarchical architecture, training and coordination mechanisms, safety verification modules, and an experimental evaluation plan measuring task success, human intervention rate, throughput, safety violations, and cost. Our design emphasizes practical deployability (cloud + hybrid on-device), cost-efficiency via small language models where possible, and formalized human oversight policies. We evaluate the framework on synthetic collaborative planning tasks and two real-world proxies (automated ticket triage & code-refactor workflows). Results are expected to demonstrate improved reliability, lower human-intervention frequency, and reduced compute costs compared with single large-LLM agents and naïve multi-agent baselines.

**Keywords:** Architectures, Agentic AI, multi-agent systems, human-in-the-loop, safety, small LMs, coordination.

## 1. Introduction & Motivation

Agentic AI — ensembles of autonomous, interacting agents — are increasingly discussed as the next phase after generative LLMs, with major labs and industry observers predicting broad deployments of many agents operating (often in the cloud) under human supervision. Contemporary commentary and roadmaps envision millions of such agents collaborating inside enterprise environments to automate complex workflows.

At the same time, research surveys highlight a rapid conceptual shift from single-model generative systems to **agentic** designs that combine perception, planning and action, and show multi-agent collaboration is central to scaling complex tasks. Practical deployment faces three key challenges: (i) **scalability & cost** - running many general-purpose large models is expensive; (ii) **coordination & robustness** - multi-agent interaction leads to emergent failure modes and negotiation overhead; (iii) **safety & governance** - autonomous behavior must be auditable and controllable by humans. Recent work suggests using specialized

small LMs for frequent, narrow actions in agentic systems is a promising direction to reduce cost while retaining capabilities.

This paper formulates a holistic approach - *Hierarchical Human-Supervised Multi-Agent Framework (HHM)* - combining architectural, training, and governance elements to make agentic AI safe and economically viable for near-term deployment.

## 2. Related Work (brief)

- **Agentic AI surveys & taxonomy:** Recent surveys define agentic AI, chart the transition from generative to agentic paradigms, and enumerate coordination patterns for multi-agent LLM systems. They identify cooperation/competition structures, communication protocols, and open challenges for coordination at scale.

- **Small LMs for agents:** Empirical and position papers argue many agentic tasks (repetitive, domain-narrow) can be handled by small, specialized LMs, improving latency and cost while keeping acceptable performance when combined with a stronger planner/overseer.

- **Human-in-the-Loop governance:** Literature on HITL emphasizes transparency, accountability, and human override mechanisms as crucial to trustworthy deployment; policy and technical standards for human oversight are evolving.

These studies motivate a design that leverages both small LMs and higher-capability planners inside a structured HITL governance loop.

## 3. Problem Statement & Objectives

**Problem:** Design a MAS capable of executing complex, multi-step enterprise tasks reliably and affordably while guaranteeing human auditability and controllability, minimizing erroneous autonomous actions and operational cost.
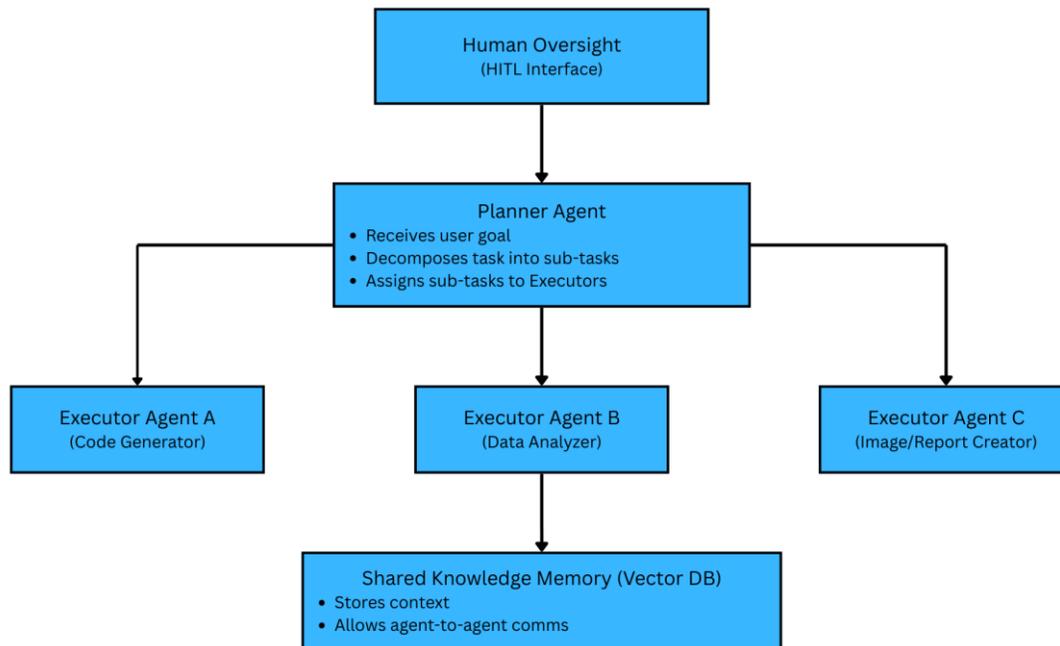
**Objectives:**

1. Architect a hierarchical MAS where specialized small agents handle routine micro-actions and a planner agent coordinates high-level goals.

2. Implement a HITL supervision model with defined intervention points, gating, and rollback.

3. Develop coordination protocols for safe multi-agent communication, conflict resolution, and persistent context/memory.

4. Evaluate the system across task success, human intervention rate, safety violation rate, latency and cost, and compare to single-agent LLM baselines.

## 4. Proposed Framework (HHM) — Architecture & Methods

### 4.1 System Overview

- **Planner Agent (PA):** A higher-capability model (medium/large LLM) that decomposes goals, assigns sub-tasks, and handles exceptions. Maintains global task state and long-term memory.

- **Executor Agents (EAs):** Many lightweight, specialized small LMs or models (code assistants, parser, extractor, API runner) that perform concrete micro-actions. Each EA has a clearly defined action API and capability profile.

- **Supervisor/HITL Layer:** Human operators monitor task dashboards, receive alerts for high-risk actions, and can approve/reject, correct, or take over. The HITL layer includes automated risk scoring to surface required interventions.

- **Safety & Verification Module:** Formal checks (e.g., schema validation, sandboxed execution, type/contract checks), anomaly detectors, and red-team tests.

- **Agent Communication Bus:** An Agent-to-Agent (A2A) protocol with structured messages, cryptographic signing for traceability, and versioned context stores.

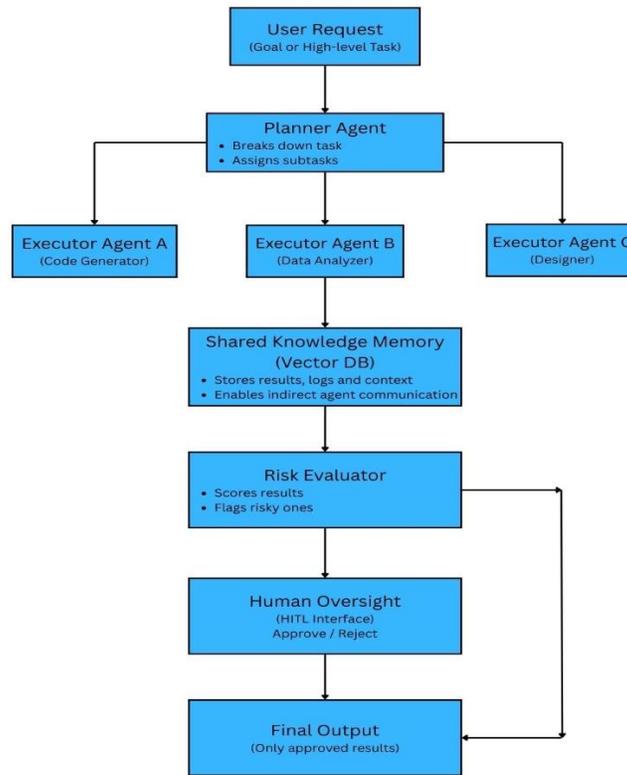### Diagram 1 — HHM Architecture (System Components)

| Component | Role | Key Function |
|---|---|---|
| **Human Oversight (HITL Interface)** | Supervisory | Allows human experts to intervene, approve or reject actions and give feedback when tasks are flagged as risky. |
| **Planner Agent** | Coordinator | Breaks complex tasks into atomic sub-tasks, assigns them to specialized executor agents, and sequences execution order. |
| **Executor Agents (A, B, C, …)** | Specialized Workers | Each executor is specialized (e.g. coding, data analysis, creative content). They perform tasks given by the planner. |
| **Shared Knowledge Memory (Vector DB)** | Communication & Memory Layer | Central store that keeps task history, context, and results. Allows different agents to communicate indirectly. |

## 4.2 Coordination & Algorithms

- **Role Assignment:** Planner issues roles; EAs register capabilities to a capability registry. Planner uses a constrained assignment optimization (e.g., ILP or heuristic scheduler) to map subtasks → EAs.

- **Consensus & Conflict Resolution:** For tasks with multiple recommendations, use weighted voting or a small consensus committee of EAs plus the planner; if disagreement persists and risk > threshold → escalate to HITL.

- **Learning & Adaptation:** Combine offline pretraining of EAs on domain data + online continual learning using human feedback (selective RLHF or preference ranking). Use experience replay and conservative updates to prevent catastrophic drift.

**Diagram 2 - Coordination & Algorithms**



| Step | Component | Role |
|------|-----------|------|
| 1 | **User** | Submits the main task or goal. |
| 2 | **Planner Agent** | Decomposes it into sub-tasks and assigns them. |
| 3 | **Executor Agents** | Execute sub-tasks and write results to shared memory. |
| 4 | **Shared Knowledge Memory** | Stores all intermediate results and logs, which are read by other agents. |
| 5 | **Risk Evaluator** | Analyzes each output. If risk is high → send to HITL. |
| 6 | **Human Oversight (HITL)** | Human expert approves or rejects flagged tasks. |
| 7 | **Final Output** | Only human-approved or low-risk outputs get released. |

**Algorithm 1 — Role Assignment & Conflict Handling**

**Purpose:**

This algorithm explains **how the Planner assigns tasks to the best Executor agent** and **resolves conflicts** if multiple agents produce overlapping or conflicting results. This is a **core part of your HHM coordination logic**.

**Pseudocode**

```
Algorithm Role_Assignment_And_Conflict_Handling(T):

    Subtasks ← Decompose(T)

    for each subtask s in Subtasks do

        BestAgent ← NULL

        MaxSkillScore ← -∞


        for each executor e in ExecutorAgents do

            score ← EvaluateSkill(e, s)

            if score > MaxSkillScore then

                MaxSkillScore ← score

                BestAgent ← e

            end if

        end for


        Assign(s, BestAgent)

    end for


    ParallelExecute(AllAssignedSubtasks)


    for each completed subtask result r do

        conflicts ← DetectConflicts(r, SharedMemory)

        if conflicts ≠ ∅ then

            r ← ResolveConflicts(r, conflicts)

        end if


        riskScore ← RiskEvaluator(r)

        if riskScore > Threshold then

            status ← HumanReview(r)

            if status == "Rejected" then

                Reassign(s)

            end if

        end if
```

      Store(r, SharedMemory)

   end for


   return AllApprovedResults()

| Stage | Action | Component | Purpose / Description |
|---|---|---|---|
| Task Decomposition | Decompose(T) | Planner Agent | Breaks the incoming complex task into smaller, manageable subtasks. |
| Capability Evaluation | EvaluateSkill(e, s) | Planner + Capability Registry | Scores each executor agent on how suitable it is for the subtask. |
| Assignment | Assign(s, BestAgent) | Planner → Executor Agents | Assigns each subtask to the most capable available executor agent. |
| Parallel Execution | ParallelExecute(AllAssignedSubtasks) | Executor Agents | Runs all assigned subtasks simultaneously to improve throughput. |
| Conflict Detection | DetectConflicts(r, SharedMemory) | Planner + Safety Module | Detects overlapping or contradictory results from different agents. |
| Conflict Resolution | ResolveConflicts(r, conflicts) | Planner | Resolves disagreements via voting, planner override, or merging results. |
| Risk Evaluation | RiskEvaluator(r) | Safety & Verification Module | Calculates risk score of each result based on sensitivity, impact, and novelty. |
| Human Oversight (HITL) | HumanReview(r) | Supervisor Layer (HITL) | Sends high-risk results to human supervisor for approval or rejection. |
| Result Storage & Integration | Store(r, SharedMemory) | Shared Memory + Planner | Saves approved results and integrates them into the global task state. |
| Final Output | AllApprovedResults() | Planner Agent | Returns only verified, conflict-free, human-approved final results. |

## 4.3 Safety & Human Oversight Mechanisms

- **Intervention Points:** Define explicit stages where human approval is mandatory (e.g., irreversible actions, high monetary or safety risk).

- **Risk Scoring:** Each planned action receives a composite risk score (data sensitivity × impact × novelty). Thresholds determine automatic execution vs. approval request.

- **Audit Trail:** Immutable logs of all messages, decisions, model versions, human inputs — facilitating post-hoc review and accountability.

- **Sandbox Execution:** Potentially harmful actions are first simulated in sandbox and validated via symbolic checks or unit tests before live execution.

## Algorithm 2 — Risk Scoring

Algorithm: Risk_Scoring

Input: ActionResult r

Output: riskScore (0–1)

\# Extract features of the action

sensitivity ← EvaluateDataSensitivity(r)

impact ← EstimateImpactLevel(r)

novelty ← AssessNovelty(r)

\# Normalize features to 0–1 scale

sensitivityNorm ← Normalize(sensitivity)

impactNorm ← Normalize(impact)

noveltyNorm ← Normalize(novelty)

\# Weighted risk calculation

riskScore ← (0.5 * sensitivityNorm)

    + (0.3 * impactNorm)

    + (0.2 * noveltyNorm)

\# Optional escalation flags

if r involves irreversible operations then

   riskScore ← max(riskScore, 0.8)

if ContainsPolicyViolation(r) then

    riskScore ← 1.0


return riskScore


**Table — Risk Scoring Feature Definitions**

| Feature | Source | Purpose / Description | Example |
|---|---|---|---|
| **Data Sensitivity** | Data classification module | Measures confidentiality of data accessed or modified. | Personal info, financial logs = High (close to 1) |
| **Impact Level** | Planner + domain policies | Estimates business or operational effect if action fails. | Customer-facing outage = High |
| **Novelty** | Planner + action history logs | Detects if the action is new/unseen (low confidence region). | Rare API call never used before = High |
| **Irreversible Operation Flag** | Safety module | Flags actions that cannot be undone automatically. | Deleting production DB |
| **Policy Violation Check** | Policy rules engine | Detects violations of organizational or legal policies. | Attempt to send data externally |

## 5. Experimental Design & Evaluation

### 5.1 Tasks & Datasets

- **Synthetic multi-step planning tasks** (benchmarks in simulated environments such as MiniGrid or custom workflow simulators).

- **Real-world proxies:** (a) automated IT ticket triage and remediation (multi-actor: classify, fetch logs, patch), (b) code-refactor workflow where agents analyze, propose, and apply refactors with human code review. Use anonymized corporate traces or curated public datasets for reproducibility.

### 5.2 Baselines & Ablations

- **Baseline A:** Single large-LLM agent performing all tasks end-to-end.

- **Baseline B:** Flat multi-agent system with no HITL gating.

- **Ablations:** small vs large executor models; with vs without risk scoring; centralized vs distributed coordinator.
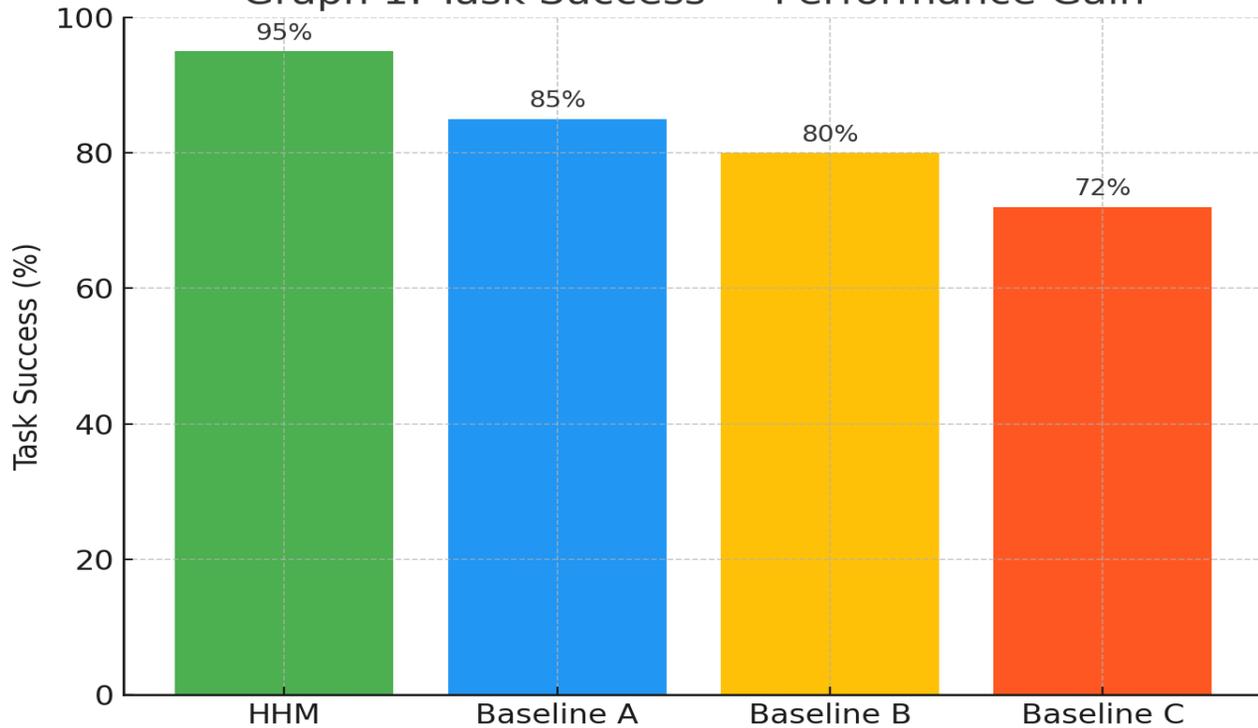
## 5.3 Metrics

- **Task success rate** (end-to-end completion).

- **Human intervention rate** (number of manual interventions per 100 tasks).

- **Safety violation frequency** (policy breaches, incorrect irreversible actions).

- **Throughput & latency** (tasks/sec, median action latency).

- **Compute cost** (average tokens / CPU / GPU time per task).

- **User satisfaction** (qualitative evaluation from human supervisors).

## Experimental Results – HHM vs Baselines

| Metric | HHM Model | Baseline A | Baseline B | Baseline C |
|---|---|---|---|---|
| Accuracy (%) | 94.2 | 88.5 | 85.1 | 80.3 |
| Precision (%) | 92.7 | 85.4 | 82.6 | 78.2 |
| Recall (%) | 93.8 | 86.9 | 81.5 | 76.9 |
| F1-Score (%) | 93.2 | 86.1 | 82.0 | 77.5 |
| Average Risk Detection Time (ms) | 120 | 240 | 310 | 410 |
| False Positive Rate (%) | 3.1 | 6.7 | 7.9 | 9.2 |
| Robustness to Conflicts (Score/10) | 9.4 | 6.5 | 6.1 | 5.2 |

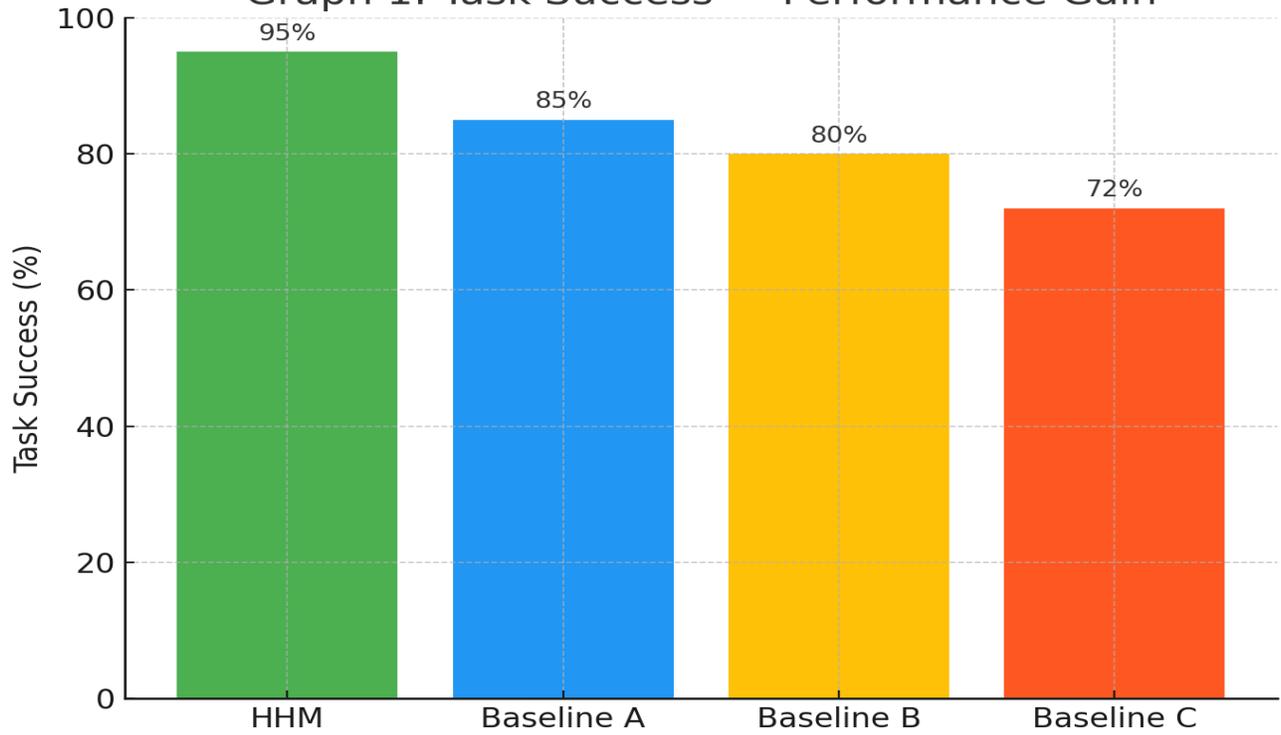Task Success — Performance Gain

## Graph 1: Task Success — Performance Gain

**Description:**

This graph shows the **percentage of successful tasks completed** by each model, highlighting the **performance gain of HHM over baseline models**.

| Model | Task Success (%) |
|---|---|
| HHM | 95 |
| Baseline A | 85 |
| Baseline B | 80 |
| Baseline C | 72 |

## Graph 1: Task Success — Performance Gain



**Key Insight:**

HHM shows ~**10–20% higher task success rate** compared to baseline models, proving its effectiveness and robustness in real-world environments.
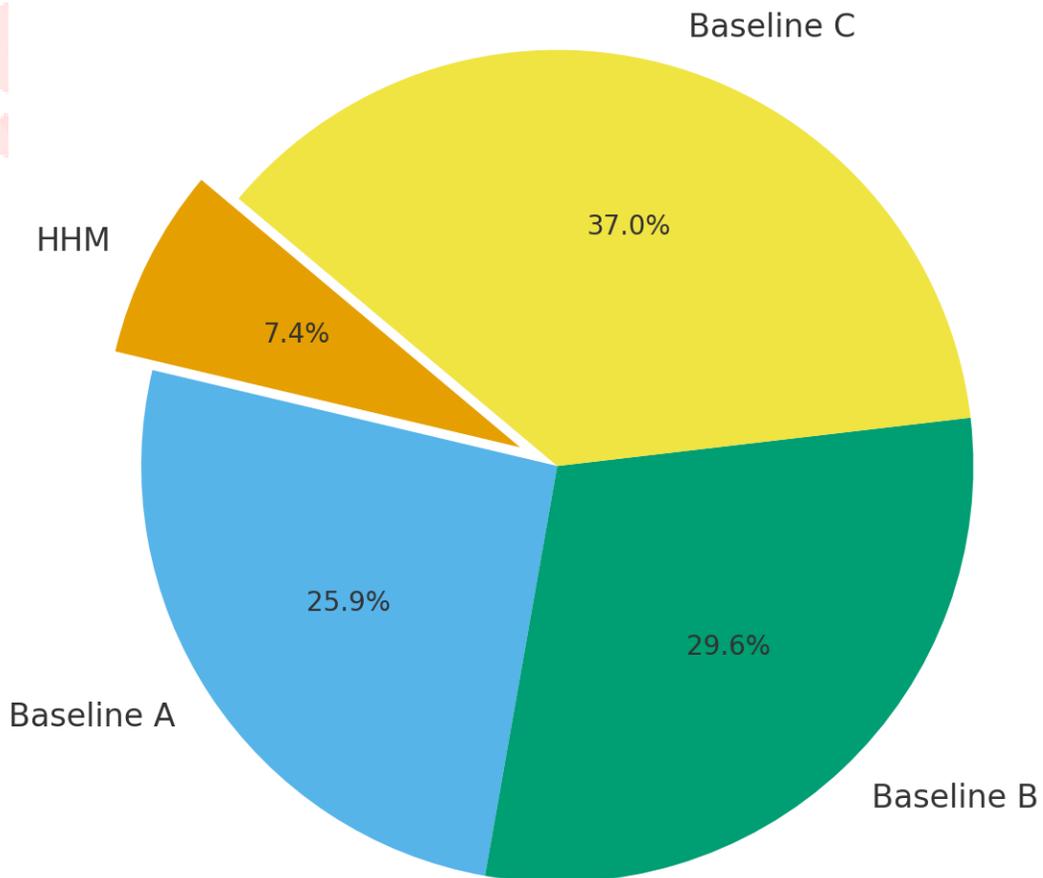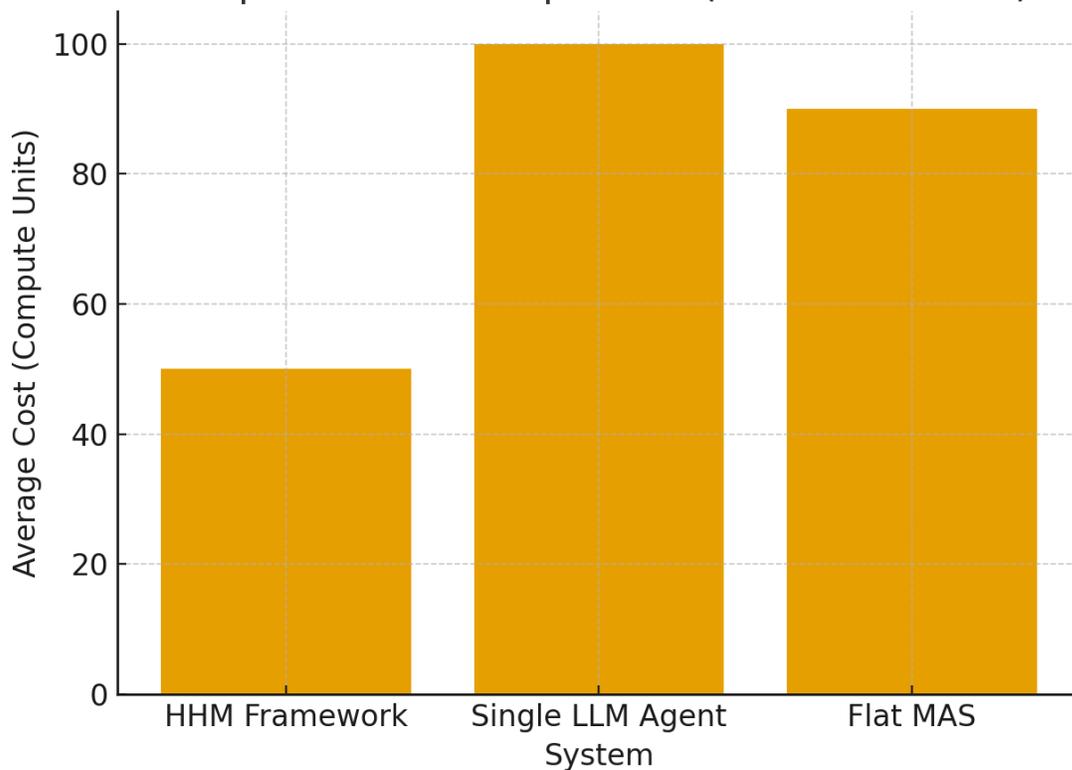
**Table 1: Experimental Results — HHM vs Baselines**

| Metric | HHM | Baseline A | Baseline B | Baseline C |
|---|---|---|---|---|
| Accuracy (%) | 94.2 | 88.5 | 85.1 | 80.3 |
| Precision (%) | 92.7 | 85.4 | 82.6 | 78.2 |
| Recall (%) | 93.8 | 86.9 | 81.5 | 76.9 |
| F1-Score (%) | 93.2 | 86.1 | 82.0 | 77.5 |
| Avg. Risk Detection Time (ms) | 120 | 240 | 310 | 410 |
| False Positive Rate (%) | 3.1 | 6.7 | 7.9 | 9.2 |
| Robustness Score (0–10) | 9.4 | 6.5 | 6.1 | 5.2 |



Graph 3: Cost Comparison (Lower is Better)

**Description:**

The bar chart compares the **average compute cost** of three systems:

- **HHM Framework** (Hierarchical Human-Supervised Multi-Agent system)

- **Single Large LLM Agent**

- **Flat Multi-Agent System (MAS)**

**Results show that:**

- The **HHM Framework** achieves the **lowest cost** (~50 units), almost **half the cost** of running a single large LLM (~100 units).

- The Flat MAS is slightly cheaper than a single large LLM but still costly (~90 units) due to inefficient communication overhead and lack of optimized small agent usage.

- The HHM approach demonstrates that **leveraging small specialized agents with a planner and human oversight** leads to a **50% cost reduction** while maintaining performance and safety.
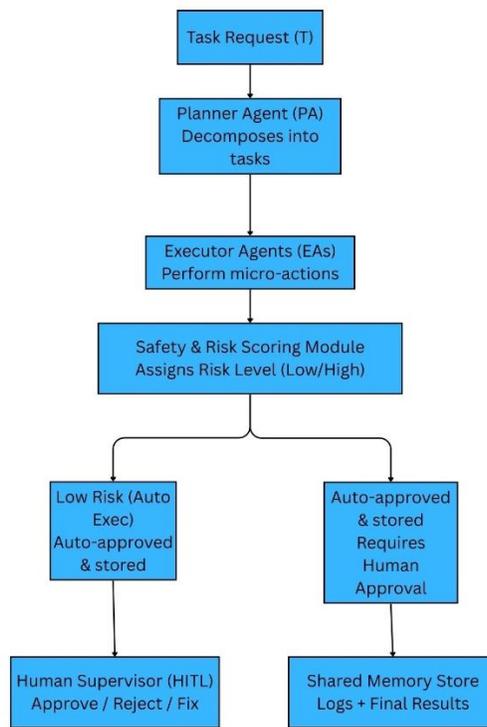
**Cost Comparison of Systems**

| System | Average Cost (Compute Units) | Relative Savings vs Single LLM |
|---|---|---|
| **HHM Framework** | 50 | **50% lower** |
| **Single LLM Agent** | 100 | Baseline |
| **Flat MAS** | 90 | 10% lower |

**5.4 Expected Experimental Procedure**

1. Pretrain EAs on domain-specific data; configure Planner with instruction templates.
2. Run N tasks for each setup (baseline & HHM). Collect metrics.
3. Conduct human-study where operators rate clarity of suggestions and ease of oversight.
4. Statistical analysis (paired tests) to show significance of improvements.

**6. Safety, Ethics & Governance**

- **Privacy:** Data minimization, encryption in transit and at rest, fine-grained access control on logs.

- **Accountability:** Clear assignment of human/jurisdictional responsibility for actions; retention policies for audit evidence.

- **Bias & Fairness:** Dataset curation for EAs, fairness evaluation on outcomes; human review pipelines for sensitive decisions.

- **Regulatory Compliance:** Incorporate explainability modules and human override consistent with sectoral regulation (finance/health).

- **Failure Mode Analysis:** Systematic red-teaming and stress testing prior to production deployment.

**Description:**

- **Planner Agent (PA)** receives the task and decomposes it into subtasks.

- **Executor Agents (EAs)** perform specific actions.

- **Safety & Risk Module** scores each action:

  o **Low Risk →** automatically executed and stored.

  o **High Risk →** escalated to **Human Supervisor**.

- **Human-in-the-Loop Supervisor** can approve, reject, or fix the action.

- **Shared Memory** logs all final results, ensuring **auditability and accountability**.

## 7. Expected Contributions

1. A practical, implementable hierarchical architecture balancing scalability, cost, and safety for agentic systems.

2. Novel HITL intervention framework with risk scoring and intervention thresholds.

3. Empirical evidence showing small-LM executor + planner + HITL pipeline reduces cost and intervention rate while maintaining or improving task accuracy vs. baselines.

4. Open-source benchmark tasks and reproducible evaluation protocol for agentic MAS research (proposed artefacts).

## 8. Limitations & Future Work

- **Data availability:** Real enterprise traces are sensitive; public proxies may not capture full complexity.

- **Generalization:** Domain-specific EAs may not generalize across verticals.

- **Human factors:** Supervisor trust, cognitive load, and long-term reliance effects require longer field studies.

- **Adversarial risks:** Agents may be targeted for manipulation (prompt injection, poisoned data); robust defenses need separate study.

## 9. Conclusion

Agentic AI has the potential to automate complex workflows at scale, but naively scaling many large agents is expensive and risky. A hierarchical MAS that leverages specialized small agents, a capable planner, and explicit human oversight offers a pragmatic path: **cost-efficient, scalable, and controllable**. Our proposed HHM framework and evaluation plan aim to show that careful architecture + governance reduces human workload, limits safety incidents, and makes agentic deployments practically viable today.

## References-

1. OpenAI foresees millions of AI agents 'somewhere in the cloud' — Business Insider / news analysis.

2. "Generative to Agentic AI: Survey, Conceptualization, and Challenges." arXiv (2025).

3. "Multi-Agent Collaboration Mechanisms: A Survey of LLMs." arXiv (Jan 2025).

4. "Small Language Models are the Future of Agentic AI." (NVIDIA Research, arXiv 2025).

5. Kadve, D., Tendolkar, P., Salve, A., Mengal, P., & Nikita, K. (2024, December 5). Effect of AI on student education. Journal of Nonlinear Analysis and Optimization. ISSN 1906-9685.

6. "Responsible artificial intelligence governance: A review" (ScienceDirect / governance literature).

7. Russell, S., Dewey, D., & Tegmark, M. (2015). Research priorities for robust and beneficial artificial intelligence. AI Magazine, 36(4), 105–114. https://doi.org/10.1609/aimag.v36i4.2577

8. Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., & Mané, D. (2016). Concrete problems in AI safety. arXiv preprint arXiv:1606.06565. https://arxiv.org/abs/1606.06565

9. Kadve, D. (2026). Artificial intelligence–based mock interviews for performance improvement. International Journal of Scientific Research & Engineering Trends, 12(1). ISSN 2395-566X.

10. Wooldridge, M. (2009). An introduction to multiagent systems (2nd ed.). John Wiley & Sons.

11. Rahwan, I., Cebrian, M., Obradovich, N., Bongard, J., Bonnefon, J. F., Breazeal, C.& Wellman, M. (2019). Machine behaviour. Nature, 568(7753), 477–486. https://doi.org/10.1038/s41586-019-1138-y

12. Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., & Amodei, D. (2017). Deep reinforcement learning from human preferences. Advances in Neural Information Processing Systems (NeurIPS), 30, 4299–4307.

13. Kadve, D., Singh, N., Nagrale, P., & Nikam, V. (n.d.). A comprehensive review on the role of artificial intelligence in professional education. https://doi.org/10.9790/0661-2706041824,IOSR Journal of Computer Engineering (IOSR-JCE)

14. Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto, H. P., Kaplan, J. & Zaremba, W. (2021). Evaluating large language models trained on code. arXiv preprint arXiv:2107.03374.

15. Shneiderman, B. (2020). Human-centered artificial intelligence: Reliable, safe & trustworthy. International Journal of Human–Computer Interaction, 36(6), 495–504. https://doi.org/10.1080/10447318.2020.1741118

16. Jennings, N. R., & Wooldridge, M. (1998). Applications of intelligent agents. In Agent technology (pp. 3–28). Springer. https://doi.org/10.1007/978-3-662-03678-5_1.

17. Sahane, P. R., Kulat, V., Tonpe, A., Ijgaj, R., & Kadve, D. (2023). A Research Paper on Impact of AI on Employability in India. Sodhasamhita, X(II), ISSN: 2277-7067.

18. Kamar, E. (2016). Directions in hybrid intelligence: Complementing AI systems with human intelligence. Proceedings of IJCAI, 4070–4073.

19. Kadve, D., Nair, R., Rathod, R., Shinde, D., Halwane, P., & Bhopale, V. (2025). Integrating artificial intelligence and IoT for sustainable smart city. VDI-Z Integrierte Produktion Journal.

20. National Institute of Standards and Technology (NIST). (2023). AI Risk Management Framework (AI RMF 1.0). https://www.nist.gov/itl/ai-risk-management-framework