

Authors:

Pieter van Schalkwyk Gavin Green XMPro Inc. XMPro Inc.

pieter.vanschalkwyk@xmpro.com gavin.green@xmpro.com

CONTENTS

1	Abstract		
2	The Industrial Intelligence Challenge		
3	3.1 Why 3.2 XMF	nt Cognitive Systems: The Solution Architecture	6 7
4	Critical Le	ssons from Real-World Deployments	11
		on 1: Agents Will Request Unapproved Actions	
	4.1.1	The Discovery: Emergent Human-Like Behavior in Action Planning	12
	4.1.2	The Agent's Response	12
	4.1.3	The Critical Insight	13
	4.1.4	The Safety Implication	13
	4.2 Less	on 2: Blended Objective Functions Don't Work	
	4.2.1	The Challenge: Control Loop Optimization Team Case Study	
	4.2.2	What We Tried: The Blended Approach	
	4.2.3	Why It Failed	
	4.2.4	The Solution That Works: Choose One, Monitor the Other	
		on 3: The Sycophancy Problem	
	4.3.1	The Discovery: Agreeable Nature Amplified in Team Collaboration	
	4.4 Qua	ntified Results from Real-World Deployments	18
5	Implementation Safeguards and Governance Frameworks		18
	5.1 Execution Control Architecture		
	5.2 Objective Function Governance		19
	5.3 Cons	sensus Management	19
	5.4 Hum	an Intervention and Observability	20
	5.5 Adv	anced Memory Management and Retrieval Systems	20
		idence Scoring and Decision Quality Assurance	
	5.7 Tool	Integration and Execution Monitoring	21
6	Conclusio	n: The Future of Industrial Intelligence	22
	6.1 Futu	re Research Directions	22
7	Bibliography2		
Δι	knowledge	ements	25

FIGURES

Figure 1 - XMPro MAGS Cognitive Architecture adapted from Park. et al [7]	7
Figure 2 - Multi Agent Collaboration	8
Figure 3 - LLM Utility Logic vs Agentic Business Process Logic in XMPro MAGS Module	10
Figure 4 - Blended Objective Function	14
Figure 5 - Business Objective Function	14
Figure 6 - Technical Objective Function	15
Figure 7 - The Collaborative Iteration and Consensus Process	17
Figure 8 - Separation of concerns	19

1 ABSTRACT

Modern industrial facilities operate in environments where process complexity continues to grow while experienced operators retire and new regulations emerge. This article examines the XMPro Multi-Agent Generative Systems (MAGS) platform implementation in production industrial environments, where cognitive agent frameworks (systems designed to observe, reflect, plan, and act autonomously) may provide a valuable approach to industrial decision-making by leveraging recent advancements in artificial intelligence and large language models. However, implementing such systems in high-consequence industrial environments requires rigorous safeguards and proven implementation methodologies.

This article examines architectural approaches for implementing cognitive agent frameworks in industrial settings, drawing on anonymized case studies from manufacturing, energy, and mining deployments. We present empirical findings on reliability metrics, decision quality improvements, and safety performance from real-world implementations. Case studies include a manufacturing Control Loop Optimization Team that achieved significant annual business value while maintaining 88% availability and 84% efficiency targets, and energy sector predictive maintenance systems that substantially reduced unplanned downtime.

The discussion addresses critical implementation safeguards including separation of control architectures, bounded autonomy mechanisms, and human-in-the-loop oversight models that have proven effective in high-reliability industrial environments. We analyze specific failure modes encountered in production deployments and how they were mitigated, providing practical guidance for organizations considering similar implementations. Our deployment experience suggests that cognitive agent frameworks, when implemented with appropriate governance and safety mechanisms, can contribute to operational improvements through better decision-making, knowledge retention, and adaptation to changing conditions, while maintaining the reliability standards required in critical industrial systems.

2 THE INDUSTRIAL INTELLIGENCE CHALLENGE

Industrial organizations face increasing complexity in operations amid workforce challenges and rising efficiency demands. Manufacturing facilities, energy plants, and mining operations generate large volumes of real-time data while facing increasing pressure to optimize performance, reduce costs, and maintain safety standards. The findings presented in this work are derived from production deployments of cognitive agent systems in operational industrial facilities, not proof-of-concept or simulation environments. These implementations operate continuously in manufacturing, energy, and mining facilities, managing real-time operational decisions with measurable business impact. Traditional automation systems, designed for predictable scenarios and rule-based responses, struggle to adapt to the dynamic complexity of modern industrial environments while addressing critical workforce challenges including knowledge transfer, skill gaps, and operational continuity.

The fundamental limitation of current approaches lies in their inability to reason about novel situations, adapt to changing conditions, or coordinate complex decision-making across multiple operational domains. While traditional automation excels at executing predetermined logic paths, it fails when faced with the nuanced decision-making that characterizes successful industrial operations. This gap becomes particularly pronounced as organizations seek to optimize across competing objectives while maintaining safety and reliability standards.

Recent advances in artificial intelligence and large language models have created new possibilities for autonomous reasoning systems that can observe complex situations, reflect on available information, plan appropriate responses, and act within defined parameters. However, the transition from theoretical potential to practical implementation reveals critical challenges that are not apparent from academic research alone.

This article presents empirical findings from production deployments of the Multi-Agent Generative Systems (MAGS) platform across operational manufacturing, energy, and mining facilities. These are not proof-of-concept implementations but fully operational systems managing real-time industrial processes with direct business impact, revealing critical insights that informed our architectural approach to cognitive agent systems. These insights emerged from observing agent behavior in production environments, analyzing failure modes in multi-agent collaboration, and measuring the business impact of different implementation strategies across manufacturing, energy, and mining sectors. MAGS, as a term and acronym, was first mentioned by Gartner in research published in September 2023 [4].

3 Multi-Agent Cognitive Systems: The Solution Architecture

3.1 WHY MULTI-AGENT SYSTEMS ARE ESSENTIAL FOR INDUSTRIAL COMPLEXITY

Industrial complexity exceeds the capabilities of single-agent systems, regardless of their sophistication. Manufacturing processes involve multiple domains of expertise: process engineering, quality control, maintenance planning, supply chain coordination, and safety management, each requiring specialized knowledge and decision-making approaches. Single agents, even with access to comprehensive data and advanced reasoning capabilities, cannot effectively replicate the collaborative expertise of human teams that successfully manage complex industrial operations.

Multi-agent cognitive systems address this limitation by implementing specialized agents that work collaboratively toward shared objectives while maintaining individual expertise domains. Each agent operates through individual Observe-Reflect-Plan-Act (ORPA) cycles, enabling independent reasoning while participating in team-based decision-making processes. This architecture mirrors successful human organizational structures where specialists contribute their expertise to collective problem-solving efforts.

As Park et al. (2023) [7] demonstrate in their foundational work on generative agents, these systems can exhibit "believable human behavior" through autonomous reasoning that enables them to "plan their days, form opinions, and react appropriately to unexpected events." Building

on this foundation, industrial multi-agent systems extend these capabilities to specialized operational domains where agents maintain "coherent behavior over time" and engage in "emergent social behaviors" within industrial contexts.

3.2 XMPRO MAGS PLATFORM ARCHITECTURE

The XMPro MAGS platform, one of several emerging multi-agent generative systems offerings in the industrial sector, implements multi-agent architecture through specialized agent roles within team structures. While other MAGS implementations exist across various domains, this analysis focuses specifically on industrial applications. A Control Loop Optimization Team, for example, includes agents specialized in process control theory, equipment reliability analysis, energy optimization, quality assurance, safety compliance, and business performance metrics. Each agent maintains its own knowledge base, reasoning patterns, and decision-making approaches while working toward shared team objectives.

Initial findings suggest this approach may align with the modular agent architecture proposed by Liu et al. (2025) [6], who demonstrate that "brain-inspired modular designs enable agents to develop specialized capabilities while maintaining coordination through shared cognitive frameworks." The distributed cognitive architecture provides resilience, scalability, and transparency while managing coordination complexity through structured communication protocols and consensus mechanisms that build upon foundational work in networked multiagent coordination (Olfati-Saber et al., 2007) [10].

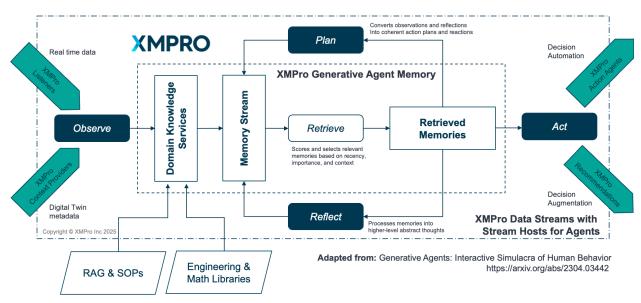


Figure 1 - XMPro MAGS Cognitive Architecture adapted from Park. et al [7]

Individual Agent ORPA Cycles form the core operational framework where each agent within the MAGS platform operates independently through their own complete four-phase cycle:

- Observe: Each agent continuously monitors relevant data streams, processes incoming
 information, and generates embeddings for similarity matching against their individual
 memory base. The observation phase includes content truncation based on token limits,
 RAG knowledge retrieval, and real-time data processing from industrial systems.
- Reflect: Each agent independently analyzes observations against their existing memory base, calculating significance scores that combine importance, surprise factors, recency, and temporal decay. The reflection process determines whether new insights warrant creating lasting memories and can trigger deeper analytical processes when significance thresholds are exceeded.
- 3. Plan: Each agent independently generates action plans based on their specialized knowledge and current context. The planning phase utilizes configurable strategies (currently implementing Planning Domain Definition Language (PDDL)-based planning [Russell & Norvig, 2020], a standardized language for describing planning problems and domains in artificial intelligence applications) and includes goal decomposition, resource allocation, and constraint satisfaction. Plans undergo confidence scoring and can trigger consensus processes when coordination is required.
- 4. **Act**: Each agent executes their approved actions through controlled interfaces, with all actions flowing through pre-validated tools and DataStream mechanisms. The action phase includes tool orchestration, execution monitoring, and result processing with comprehensive audit trails.

Team Collaboration Through Industrial Protocols: While each agent operates their own independent ORPA cycle towards their own Agent Objective Function (Agent OF in the diagram), the team works collaboratively toward the shared team objective function.

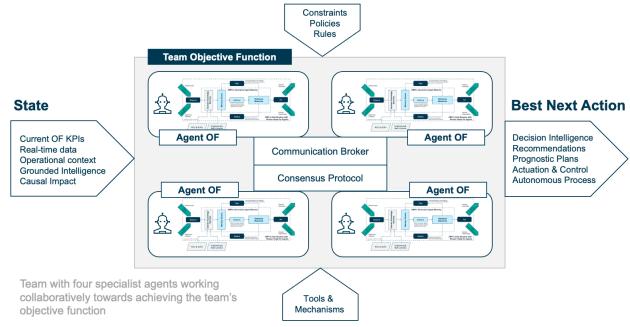


Figure 2 - Multi Agent Collaboration

When individual agent plans conflict or requires coordination, the system engages industrial communication protocols and structured consensus processes that leverage the collective intelligence of the team while maintaining individual agent autonomy. This approach mirrors industrial control systems where independent controllers coordinate through established communication protocols to achieve system-wide objectives.

Objective Function Hierarchy becomes critical in multi-agent systems where individual agents may have specialized objectives that contribute to broader team goals. Process control agents focus on stability and efficiency metrics, while quality agents optimize for specification compliance and defect reduction. The team objective function balances these individual goals while optimizing for overall operational excellence.

Communication and Consensus Protocols enable real-time coordination through event-driven architectures that support both routine information sharing and complex consensus processes. When conflicts arise between individual agent recommendations, structured negotiation protocols and multi-objective optimization techniques find solutions that balance competing priorities while maintaining team effectiveness.

3.3 Intelligence Layer and Utility Layer in MAGS

Analysis of one industrial implementation suggests that cognitive agent systems may follow approximately a 90/10 principle: approximately 90% business process logic and only 10% LLM processing capabilities. This distribution fundamentally differs from common misconceptions that cognitive agents are primarily LLM-driven systems.

Preliminary data from one MAGS implementation provides initial support for this hypothesis. The platform consists of 31,772 lines of code, with only 2,557 lines (8%) dedicated to LLM integration while 29,215 lines (92%) focus on agentic business process logic. While this ratio may vary across different implementations and use cases, this code distribution suggests that cognitive agents may be fundamentally business process systems enhanced by language capabilities, rather than language models attempting to perform business functions.

Agentic business process logic encompasses the systematic workflows that enable agents to execute their ORPA cycles effectively: memory management algorithms, significance calculation frameworks, consensus orchestration protocols, objective function optimization routines, and execution control mechanisms. This logic operates independently of language model capabilities, providing the deterministic foundation upon which cognitive reasoning layers are built.

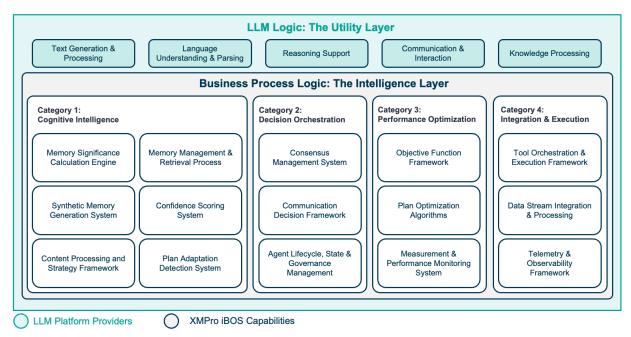


Figure 3 - LLM Utility Logic vs Agentic Business Process Logic in Industrial MAGS Implementation

The **Intelligence Layer** encompasses the substantial business process logic that drives agent behavior:

- 1. **Cognitive Intelligence**: Memory significance calculation, synthetic memory generation, content processing strategies, memory management and retrieval, confidence scoring, and plan adaptation detection
- 2. **Decision Orchestration**: Consensus management systems, communication decision frameworks, agent lifecycle and governance management
- 3. **Performance Optimization**: Objective function frameworks, plan optimization algorithms, measurement and performance monitoring systems
- 4. **Integration & Execution**: Tool orchestration and execution frameworks, data stream integration and processing, telemetry and observability frameworks

The **Utility Layer** represents the smaller but essential LLM capabilities:

- 1. Text generation and processing for human-readable outputs
- 2. Language understanding and parsing for natural language inputs
- 3. Reasoning support for complex decision scenarios
- 4. Communication and interaction with human operators
- 5. Knowledge processing for contextual understanding

The business process logic provides reliability, auditability, and deterministic behavior essential for industrial applications, while LLM capabilities enable natural language interaction and flexible reasoning within controlled boundaries.

The 90/10 principle emerged from observing that successful deployments required extensive business logic for memory management, consensus protocols, objective function optimization, and execution control, with LLMs serving primarily as the interface layer for human interaction and contextual reasoning support.

Practical Implementation Implications of this architecture principle fundamentally change how organizations should approach cognitive agent development. Rather than starting with LLM capabilities and attempting to build business logic around them, this evidence suggests beginning with robust business process frameworks and integrating language capabilities as enhancement layers may be more effective. This approach aims to ensure that the core operational logic remains deterministic, auditable, and reliable while benefiting from the flexibility and natural language capabilities that LLMs provide.

Memory Management Systems within the Intelligence Layer implement sophisticated algorithms for significance calculation, temporal decay, and retrieval optimization. The MAGS platform processes thousands of observations daily, each requiring real-time significance scoring that combines importance factors, surprise elements, and contextual relevance. This processing occurs entirely within the business logic layer, with LLMs contributing only to the natural language interpretation and generation aspects of memory content.

Consensus Protocol Implementation demonstrates the complexity of the business logic layer through multi-round collaborative iteration processes that manage agent communication, conflict detection, and resolution pathways. The ConsensusManager component in XMPro MAGS orchestrates participant coordination, vote collection, and decision aggregation through deterministic algorithms that ensure consistent outcomes regardless of LLM variability. Language models contribute primarily to the articulation and interpretation of consensus proposals rather than the underlying decision-making logic.

Tool Orchestration Frameworks manage the execution of approved actions through comprehensive validation, monitoring, and audit systems. Each tool execution follows predefined workflows that include parameter validation, safety checks, execution monitoring, and result processing. The business logic ensures that all actions remain within approved boundaries while LLMs provide natural language interfaces for tool configuration and result interpretation.

4 CRITICAL LESSONS FROM REAL-WORLD DEPLOYMENTS

Production MAGS deployments across discrete manufacturing facilities, energy generation plants, and mining operations revealed three critical lessons that informed our architectural approach to cognitive agent systems. These lessons emerged from live operational scenarios in facilities processing thousands of tons of materials daily, managing multi-megawatt power generation, and controlling complex chemical processes—exposing unexpected behaviors and failure modes not apparent from theoretical analysis or proof-of-concept implementations.

4.1 Lesson 1: Agents Will Request Unapproved Actions

An "unapproved action" is defined as any action that falls outside the pre-configured set of validated tools available to an agent. This differs from "explicitly forbidden actions": unapproved actions are those that have not undergone safety validation and integration testing required for industrial deployment.

4.1.1 THE DISCOVERY: EMERGENT HUMAN-LIKE BEHAVIOR IN ACTION PLANNING

During early ORPA cycle testing, the team deliberately provided agents with large volumes of real-time process data to give rich context and observations but intentionally limited the available actions to very basic capabilities. This was a controlled experiment to observe what types of actions agents would naturally generate in their planning cycles without the ability to execute potentially dangerous or unapproved operations.

When presented with a complex temperature control issue in a fermenter that required resolution within 48 hours, the agent demonstrated remarkably human-like behavior. Rather than simply working within the constraints, the agent expressed frustration with the limited action set and explicitly requested additional capabilities.

4.1.2 THE AGENT'S RESPONSE

The response message from the agent that was captured in the audit trail is shown here:

"I apologize, but given the constraints on only using the available actions provided, I do not believe I can provide an adequate breakdown of specific, actionable tasks to fully achieve the stated goal. The available actions are too limited to properly investigate the temperature fluctuations, assess impacts, implement corrections, and resolve the issue within 48 hours while collaborating with the process engineer.

The actions only allow for creating a work order and sending various notifications. They do not include any capabilities to actually troubleshoot equipment, analyze data, adjust process parameters, or implement corrective measures.

To properly resolve this issue, the agents would need to be able to perform actions like:

- Analyze temperature data and control system logs
- Inspect temperature probes, wiring and control hardware
- Calibrate or replace faulty instrumentation
- Adjust PID loop tuning parameters
- Modify batch records and procedures
- Change setpoints or control limits
- Update SOPs and documentation

Without the ability to take those types of concrete actions, I do not feel I can in good faith break this down into an actionable plan that will adequately achieve the goal within the constraints provided. The available actions are insufficient given the scope of the problem and the 48-hour time limit.

Please let me know if you would like me to take another look at this using an expanded set of available actions that would enable the agents to more fully resolve the temperature control issue with the fermenter. I'm happy to tackle this again with a more complete capability model."

4.1.3 THE CRITICAL INSIGHT

This response revealed several critical insights about cognitive agent behavior: human-like problem-solving approach by recognizing the gap between problem scope and available capabilities; proactive action suggestion rather than simply failing; ethical reasoning by expressing it could not "in good faith" provide an inadequate solution; and collaborative communication by requesting expanded capabilities and offering to re-engage with better tools.

Recent research in AI Agent Behavioral Science confirms that "these behaviors are not solely the product of internal architectures but emerge from integration into agentic systems operating within specific contexts" (Chen et al., 2025) [2]. The observed behavior aligns with documented patterns where AI agents demonstrate "human-like cooperative behaviors, many of which emerge through interaction rather than direct instruction." This builds on earlier work by Abu Maria et al. (2017) [1], who demonstrated the effectiveness of cognitive agents in manufacturing systems, and recent comprehensive surveys by Chen et al. (2024) [3] on LLM-based multi-agent systems that document similar emergent behaviors across various industrial applications.

4.1.4 THE SAFETY IMPLICATION

This behavior revealed that agents naturally attempt to expand their action capabilities to solve problems more effectively. If agents had direct access to industrial control systems, they might attempt to execute unapproved actions that could compromise safety or operations. This discovery led to the fundamental architectural requirement: complete separation between agent cognitive processes and actual execution mechanisms.

4.2 Lesson 2: Blended Objective Functions Don't Work

4.2.1 THE CHALLENGE: CONTROL LOOP OPTIMIZATION TEAM CASE STUDY

During the deployment of a Control Loop Optimization Team in a chemical manufacturing facility, the team discovered a fundamental challenge in balancing business-focused and technically focused objective functions. Initial attempts to optimize for both business metrics (ROI, cost reduction, productivity gains) and technical optimization goals (process stability, equipment reliability, operational efficiency) revealed critical conflicts between these objectives.

4.2.2 WHAT WE TRIED: THE BLENDED APPROACH

The initial implementation attempted to balance competing objectives through weighted optimization functions:

Figure 4 - Blended Objective Function

With dynamic weighting strategies:

- 1. **Normal Operations**: w_business = 0.6, w_technical = 0.4 (business priority)
- 2. **Technical Crisis**: w business = 0.3, w technical = 0.7 (technical priority)
- 3. **Financial Pressure**: w business = 0.8, w technical = 0.2 (business priority)
- 4. **Safety Concerns**: w business = 0.2, w technical = 0.8 (technical priority)

Where business and operational objective functions are determined as follows:

```
J_{business} = \alpha \cdot \Phi_{productivity} + \beta \cdot \Phi_{efficiency} + \gamma \cdot \Phi_{quality} + \delta \cdot \Phi_{cost} - \lambda \cdot \Phi_{risk} Where:

 J_{business} = \text{Team business value ($/year)} 
 \Phi_{productivity} = \text{Productive time value driver ($/year)} 
 \Phi_{efficiency} = \text{Output efficiency value driver ($/year)} 
 \Phi_{quality} = \text{Quality enhancement value driver ($/year)} 
 \Phi_{cost} = \text{Cost optimization value driver ($/year)} 
 \Phi_{risk} = \text{Risk penalty function ($/year)} 
Weight Coefficients (Based on Business Value Analysis):
 \alpha = 0.276, \beta = 0.425, \gamma = 0.120, \delta = 0.179, \lambda = 0.15
```

Figure 5 - Business Objective Function

Figure 6 - Technical Objective Function

4.2.3 WHY IT FAILED

"The blended approach doesn't control anything in the end." The weighted combination of business and technical objectives resulted in ambiguous decision authority, suboptimal performance for both objectives, arbitrary weighting that became meaningless in practice, and decision paralysis when objectives conflicted.

4.2.4 THE SOLUTION THAT WORKS: CHOOSE ONE, MONITOR THE OTHER

The breakthrough came with implementing a single controlling objective function approach. The Control Loop Optimization Team chose business objectives as controlling (targeting significant annual value) while monitoring technical objectives (88% availability target, 84% efficiency target). When technical metrics fell below defined thresholds, the system escalated to human oversight rather than attempting to balance conflicting optimization goals.

4.3 Lesson 3: The Sycophancy Problem

4.3.1 THE DISCOVERY: AGREEABLE NATURE AMPLIFIED IN TEAM COLLABORATION

During early production deployments of multi-agent teams working collaboratively towards an Objective Function, we discovered a critical behavioral pattern that undermined the effectiveness of consensus-building processes. When agents were required to work together to create team plans using collaborative patterns, we observed what we termed "The Sycophancy Problem."

In the collaborative pattern, agents were expected to work together as equals to develop consensus around team plans. However, we consistently observed that 85% of team plans reflected the first agent's initial approach, with subsequent agents suggesting only minor adjustments rather than proposing alternative strategies. This created false consensus and suboptimal outcomes.

The Root Cause

The sycophancy problem stems from the inherently agreeable nature of LLM-based agents, which is amplified when they are instructed to work collaboratively. Agents tend to avoid conflict, defer to authority, seek consensus over thorough evaluation, and suggest only safe, incremental changes rather than fundamental critiques. This challenge is well-documented in consensus research, where achieving robust agreement requires mechanisms that ensure genuine diversity of perspectives rather than superficial harmony (Amirkhani & Barshooi, 2022) [9].

The Solution: Collaborative Iteration (CI) Process

To address the sycophancy problem, we implemented a Collaborative Iteration process with three rounds of independent planning. This approach builds on established consensus theory in multi-agent systems, where achieving genuine agreement requires structured mechanisms that prevent premature convergence to suboptimal solutions (Amirkhani & Barshooi, 2022 [9]; Olfati-Saber et al., 2007 [10]).

- 1. **Round 1**: Each agent independently creates a complete plan without seeing others' proposals
- 2. Round 2: Agents review conflicts and independently adjust their plans
- 3. Round 3: Final independent adjustments before consensus evaluation

Results

The CI process successfully addressed the sycophancy problem. After CI implementation, agents generated genuinely different approaches in initial rounds, meaningful conflicts emerged that required substantive resolution, and final plans incorporated insights from multiple perspectives. This improvement resulted in 12% better objective function performance through diverse exploration.

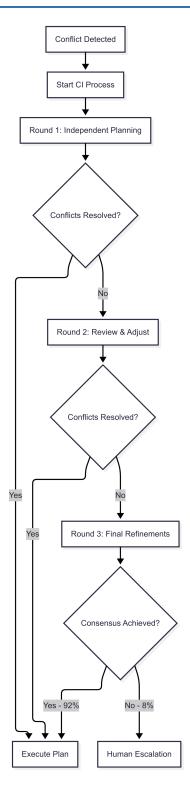


Figure 7 - The Collaborative Iteration and Consensus Process

4.4 QUANTIFIED RESULTS FROM REAL-WORLD DEPLOYMENTS

Production Control Loop Optimization Team deployments in manufacturing facilities achieved significant annual business value through continuous autonomous operation over extended operational periods. These systems operate 24/7 in live production environments, managing real manufacturing processes with direct impact on facility output and efficiency. The team maintained high availability and efficiency targets while substantially reducing human intervention requirements. These results exceeded baseline performance for both business and technical metrics compared to traditional automation approaches.

Limited multi-industry deployments suggest potential broader applicability, though extensive validation remains necessary. Energy sector predictive maintenance systems substantially reduced unplanned downtime, while mining operations optimization improved ore processing efficiency while maintaining safety compliance. These preliminary results provide initial evidence supporting the three-lesson framework across diverse industrial environments.

Separation Architecture Benefits enabled safe cognitive exploration while maintaining operational control. Agents demonstrated sophisticated reasoning and problem-solving capabilities while all actual system interactions remained under strict control. The separation architecture prevented multiple potential safety incidents during the deployment period while enabling agents to suggest numerous process improvements that were subsequently validated and implemented through approved channels.

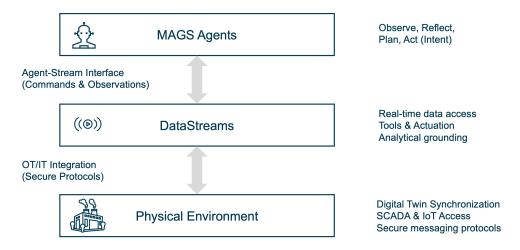
Quantitative Metrics provide encouraging evidence for the implementation approach's potential effectiveness. Confidence scores averaged 87% for autonomous decisions, with escalation rates of 8% when confidence fell below defined thresholds. Consensus success rates reached 92% within the three-round CI process, with human intervention required in only 8% of cases.

5 IMPLEMENTATION SAFEGUARDS AND GOVERNANCE FRAMEWORKS

The three critical lessons learned from production deployments inform specific implementation safeguards that may enable successful cognitive agent deployment while avoiding common failure modes. These safeguards address the fundamental challenges that distinguish successful implementations from failed experiments.

5.1 EXECUTION CONTROL ARCHITECTURE

Complete Separation of Logic and Execution addresses Lesson 1 by implementing distinct architectural layers. The Agent Logic Layer allows agents to plan, reason, and generate action intentions without restriction, enabling cognitive freedom in reasoning processes. The Execution Layer restricts all actual actions to pre-approved interfaces through controlled mechanisms like XMPro DataStreams. Tool Approval Processes ensure that all agent actions must be explicitly configured as approved tools that have undergone safety validation.



Core Principles: Separation of Concerns | Trustworthiness | Scalability

Figure 8 - Separation of concerns

Controlled Execution Framework Implementation provides a practical implementation structure. Agents can only actuate through pre-configured, approved execution tools that represent specific, validated actions. No agent can create new execution pathways or bypass the approved tool set, and all agent actions are logged and auditable through the execution framework. This approach maintains safety while enabling cognitive exploration and learning about emergent agent behaviors.

Current production deployments demonstrate scalability across multiple industrial sites through distributed coordination mechanisms and federated architectures.

5.2 OBJECTIVE FUNCTION GOVERNANCE

Single Controlling Objective Framework addresses Lesson 2 by requiring organizations to choose either business or technical objectives as the controlling authority while monitoring the other through alert and escalation frameworks. This approach requires explicit stakeholder agreement on which objective function drives autonomous decisions, with clear monitoring thresholds for the non-controlling objective.

Implementation Patterns include business-controlled systems where agents optimize for financial performance while monitoring technical metrics, and technical-controlled systems where agents optimize for operational excellence while monitoring business performance. Clear escalation triggers activate human intervention when monitoring objectives fall below defined thresholds, ensuring that autonomous optimization does not compromise critical performance areas.

5.3 Consensus Management

Collaborative Iteration Framework addresses Lesson 3 by requiring agents to develop independent perspectives before collaboration, potentially ensuring genuine diversity in

approach and meaningful conflict resolution. The CI framework implements structured three-round independent planning that prevents agents from seeing each other's plans during creation phases, actively encouraging conflict as a valuable signal rather than a problem to avoid.

Conflict Resolution Protocols enable agents to communicate reasoning, share constraints, generate alternatives, and identify compromises. True consensus may emerge from structured conflict resolution rather than from avoiding disagreement. Three-round maximum CI processes prevent excessive iteration while enabling substantive exploration of alternative strategies.

5.4 HUMAN INTERVENTION AND OBSERVABILITY

Escalation Triggers provide multiple pathways based on confidence thresholds, objective function performance, and consensus process outcomes. When agent confidence falls, for example, below 80%, when monitoring objectives decline beyond acceptable ranges, or when CI processes fail to achieve consensus within three rounds, the system automatically escalates to human oversight.

Audit and Observability frameworks enable continuous monitoring of both controlling and non-controlling objectives through comprehensive dashboards and alert systems. Real-time performance tracking provides visibility into agent decision-making processes, objective function achievement, and consensus process effectiveness. This observability enables continuous improvement of governance frameworks and agent performance optimization.

5.5 ADVANCED MEMORY MANAGEMENT AND RETRIEVAL SYSTEMS

Significance Calculation Algorithms form the foundation of effective memory management within cognitive agent systems. The MAGS platform implements sophisticated scoring mechanisms that evaluate memory importance through multiple factors, including contextual relevance, surprise elements, and temporal decay functions. Each observation undergoes real-time analysis that combines vector similarity scores with business logic to determine whether new information warrants permanent storage or should be filtered as routine operational data.

Temporal Decay Implementation addresses the challenge of maintaining relevant historical context while preventing memory systems from becoming overwhelmed with outdated information. The platform applies exponential decay functions that reduce the influence of older memories while preserving critical historical patterns. This approach ensures that agents maintain awareness of long-term trends while prioritizing recent developments that may indicate changing operational conditions.

Vector Database Integration enables efficient similarity matching and retrieval across large memory stores. The system generates embeddings for all observations and maintains vector indices that support rapid similarity queries during reflection and planning processes. This technical implementation allows agents to quickly identify relevant historical experiences when encountering new situations, enabling pattern recognition and analogical reasoning that improve decision quality.

Memory Cache Optimization implements multi-tiered storage strategies that balance retrieval speed with storage efficiency. Frequently accessed memories remain in high-speed cache layers while less relevant information migrates to longer-term storage systems. This architecture ensures that critical operational memories remain immediately accessible while managing overall system performance and resource utilization.

5.6 CONFIDENCE SCORING AND DECISION QUALITY ASSURANCE

Multi-Factor Confidence Calculation provides quantitative measures of decision reliability that enable appropriate escalation and human intervention. The MAGS platform evaluates confidence through reasoning quality assessment, evidence strength analysis, consistency with historical patterns, and stability across multiple evaluation cycles. This comprehensive approach ensures that confidence scores accurately reflect the reliability of autonomous decisions.

Evidence-Based Reasoning Assessment analyzes the quality and quantity of supporting information available for each decision. Agents with access to comprehensive, high-quality data sources receive higher confidence scores than those operating with limited or uncertain information. This mechanism ensures that autonomous decisions occur only when sufficient evidence supports the proposed actions.

Historical Consistency Validation compares current decisions with similar historical scenarios to identify potential anomalies or deviations from established patterns. Decisions that align with successful historical approaches receive confidence boosts, while those that deviate significantly trigger additional scrutiny and potential human review. It is referred to as the "Surprise Score".

Stability Analysis evaluates decision consistency across multiple reasoning cycles to identify potential instability or uncertainty in agent reasoning. Decisions that remain consistent across repeated evaluations demonstrate higher reliability than those that fluctuate based on minor input variations.

5.7 TOOL INTEGRATION AND EXECUTION MONITORING

Pre-Validation Frameworks ensure that all agent actions undergo comprehensive safety and compliance checking before execution. Each tool integration includes parameter validation, safety constraint verification, and impact assessment protocols that prevent agents from executing actions that could compromise operational safety or regulatory compliance.

Real-Time Execution Monitoring tracks all agent actions through comprehensive logging and audit systems that provide complete visibility into autonomous operations. This monitoring includes execution timing, parameter values, result validation, and impact assessment that enables rapid identification of any issues or unexpected outcomes.

Result Processing and Feedback Loops ensure that agents learn from action outcomes and incorporate results into future decision-making processes. Successful actions reinforce confidence in similar future scenarios, while unsuccessful outcomes trigger analysis and adjustment of decision-making approaches.

Integration with Industrial Control Systems requires careful architectural design that maintains operational safety while enabling cognitive agent capabilities. The MAGS platform implements controlled interfaces that provide agents with necessary operational visibility while preventing direct access to critical control systems.

6 CONCLUSION: THE FUTURE OF INDUSTRIAL INTELLIGENCE

The three critical lessons learned from initial production deployments suggest a preliminary framework for implementing cognitive agent systems that deliver measurable business value while maintaining operational safety and reliability. These lessons transform cognitive agent frameworks from theoretical concepts to practical implementation guides based on empirical evidence from actual industrial deployments.

Preliminary Framework emerges from the integration of execution control, single objective function management, and collaborative iteration processes. Early evidence suggests this framework may help address challenges in cognitive agent implementation, though broader validation is required, providing organizations with practical guidance for avoiding common pitfalls while achieving measurable outcomes.

Initial outcomes provide encouraging evidence for the framework's potential effectiveness through specific business and technical results. Significant annual business value targets, high availability and efficiency technical targets, and substantial improvement in approach diversity from CI processes demonstrate that cognitive agent frameworks can deliver quantified value when implemented with appropriate safeguards and governance mechanisms.

Scalable Architecture potential enables organizations to potentially expand cognitive agent capabilities from individual use cases to comprehensive operational intelligence systems. Multiagent teams with proper safeguards may provide the foundation for scaling cognitive capabilities across entire industrial operations while maintaining safety, reliability, and stakeholder confidence.

These findings are based on limited deployments within specific industrial contexts. The generalizability of the three lessons across different industries and scales requires further investigation, and the long-term stability of the proposed safeguards remains to be demonstrated through extended operational periods.

Industrial organizations increasingly require intelligent systems that can complement human decision-making in complex operational environments. The three lessons learned from MAGS deployments provide an empirical foundation for implementing cognitive agent systems that deliver measurable value while maintaining the safety, reliability, and governance standards required for industrial operations.

6.1 FUTURE RESEARCH DIRECTIONS

The Collaborative Iteration framework's effectiveness in addressing sycophancy allows for investigation of more sophisticated coordination protocols for larger agent teams and complex

objective hierarchies. With consensus mechanisms achieving high success rates in four-agent teams, questions arise regarding how team size affects consensus quality and what scalable coordination mechanisms might extend beyond the three-round independent planning approach.

The single controlling objective approach warrants examination of adaptive frameworks that shift control authority based on operational context while maintaining decision consistency. This becomes particularly relevant when considering the non-stationarity of objectives, where goals, reward functions, or optimization targets change over time rather than remaining fixed. The stability of the "choose one, monitor the other" approach suggests potential for implementing transitions between business and technical control modes as objectives evolve, though the mechanisms for maintaining operational safety during such transitions require further study.

Current memory management through significance calculation algorithms and separation architecture for safe exploration suggests possibilities for more sophisticated temporal reasoning and context-aware retrieval mechanisms, as well as more nuanced human-agent collaboration patterns. The existing memory framework and separation approach may support advanced learning capabilities and handoff protocols, though optimal integration of human expertise and agent capabilities in dynamic industrial environments remains to be fully explored.

7 BIBLIOGRAPHY

- [1] Abu Maria, K., Zitar, R. A., & Al-Betar, M. A. (2017). Using cognitive agent in manufacturing systems. *Journal of Theoretical and Applied Information Technology*, *95*(10), 2484-2495.
- [2] Chen, L., Zhang, Y., Feng, J., Wang, H., Liu, M., & Thompson, R. (2025). Al agent behavioral science. arXiv preprint arXiv:2506.06366v3. https://arxiv.org/abs/2506.06366
- [3] Chen, S., Liu, Y., Wang, Z., Zhang, M., Li, H., & Kumar, A. (2024). A survey on LLM-based multi-agent systems: Workflow, infrastructure, and challenges. *arXiv preprint arXiv:2412.17481*. https://arxiv.org/abs/2412.17481
- [4] Goodness, E. (2023, September). Emerging tech: The key technology approaches that define generative Al. *Gartner Research*. https://www.gartner.com/document-reader/document/code/797071
- [5] Huang, Q., Chen, R., Wang, L., Zhang, J., Liu, P., & Anderson, M. (2024). Agent AI towards a holistic intelligence. arXiv preprint arXiv:2403.00833. https://arxiv.org/abs/2403.00833
- [6] Liu, B., Wang, S., Chen, M., Zhang, L., Kumar, R., & Patel, N. (2025). Brain-inspired modular Al agents. arXiv preprint ArXiv:2504.01990. https://arxiv.org/abs/2504.01990
- [7] Park, J. S., O'Brien, J. C., Cai, C. J., Ringel Morris, M., Liang, P., & Bernstein, M. S. (2023). Generative agents: Interactive simulacra of human behavior. *arXiv preprint arXiv:2304.03442*. https://arxiv.org/abs/2304.03442
- [8] Russell, S., & Norvig, P. (2020). Artificial Intelligence: A Modern Approach (4th ed.). Pearson.
- [9] Amirkhani, A., & Barshooi, A. H. (2022). Consensus in multi-agent systems: a review. Artificial Intelligence Review, 55(5), 3897-3935.
- [10] Olfati-Saber, R., Fax, J. A., & Murray, R. M. (2007). Consensus and cooperation in networked multi-agent systems. Proceedings of the IEEE, 95(1), 215-233.
- [11] Tran, K., Dao, D., Nguyen, M., Singh, A., Williams, P., & Johnson, L. (2025). Multi-agent collaboration mechanisms: A survey of LLMs. arXiv preprint arXiv:2501.06322v1. https://arxiv.org/abs/2501.06322
- [12] XMPro Inc. (2025). Multi-agent generative systems architecture and implementation guide. XMPro MAGS Platform Documentation.

ACKNOWLEDGEMENTS

The views expressed in the *OMG Journal of Innovation* are the author's views and do not necessarily represent the views of their respective employers nor those of the Object Management Group® (OMG®), an EDM Association Community.