(REVIEW ARTICLE)

Check for updates

# Trust Architecture for Enterprise AI Assistants: Technical Mechanisms for Transparency and Security

Prem Sai Reddy Kareti *

*Jawaharlal Nehru Technological University, India.*

## Abstract

Enterprise AI assistants have become integral components of workplace software ecosystems, yet their successful adoption hinges on establishing genuine user trust. This article presents a comprehensive technical framework for implementing trust-building mechanisms within enterprise AI systems. The foundation of this framework consists of four interconnected pillars: explicit AI identity signaling, verifiable information provenance through citation systems, sensitivity-aware data handling capabilities, and secure context preservation during multi-agent handoffs. These mechanisms require thoughtful implementation across multiple layers of the technology stack, from model design to user interface components. The technical architecture proposed addresses critical enterprise requirements for transparency, reliability, and compliance while maintaining seamless user experiences. Organizations implementing these recommendations can expect increased user confidence, broader adoption, and reduced resistance to AI integration within sensitive business processes. Future developments in this domain will likely focus on standardizing trust indicators across enterprise platforms and refining context preservation during increasingly complex multi-agent workflows.

## 1. Introduction

The integration of AI assistants into enterprise software environments represents a significant shift in how knowledge workers interact with organizational systems. These assistants now perform a range of functions from answering queries and automating routine tasks to facilitating complex workflows across multiple business applications. This transition occupies "the fine balance between helping with your job and taking it," highlighting the delicate relationship between augmentation and automation in workplace contexts [1].

### 1.1. Current Landscape of AI Assistants in Workplace Software

Enterprise AI assistants have evolved from simple chatbots to sophisticated systems integrated throughout organizational software. These assistants now handle increasingly complex tasks including information retrieval, process automation, and context-aware recommendations. The workplace software ecosystem has rapidly adapted to incorporate these capabilities, with many platforms now offering AI features as core functionality rather than add-ons. This evolution reflects a broader organizational recognition that AI can serve as more than just a technological novelty and instead function as an essential workplace tool [2].

---

* Corresponding author: Prem Sai Reddy Kareti

## 1.2. The Business Case for Trusted AI Interactions

The business case for trusted AI interactions extends beyond mere efficiency gains. Organizations implementing AI assistants seek to improve employee productivity, enhance decision-making processes, and create more intuitive interfaces to enterprise systems. However, these benefits can only be realized when users genuinely trust and consequently engage with these AI systems. Successful AI integration requires "not instinctively defending against it" but rather developing frameworks that promote appropriate reliance [2]. Organizations that establish trusted AI interactions can expect increased user adoption, more efficient workflows, and ultimately better returns on their AI investments.

## 1.3. Key Challenges in Enterprise Environments

Enterprise environments present unique challenges for establishing this trust foundation. Unlike consumer applications, enterprise AI assistants often operate within highly regulated industries, handle sensitive organizational data, and must integrate with legacy systems. Additionally, enterprise users bring professional skepticism and heightened concerns about accuracy, data privacy, and the appropriate division of responsibility between human and machine. These challenges are compounded by organizational cultures that may resist automation of knowledge work and the complex technical requirements of enterprise software environments.

## 1.4. Research Questions and Scope of the Paper

This article addresses several critical questions related to building user trust in enterprise AI assistants: What technical mechanisms can effectively signal AI identity and information provenance? How can AI systems demonstrate appropriate handling of sensitive enterprise data? What architectural approaches support secure transitions between different AI subsystems? How can these mechanisms be implemented across the technology stack from back-end services to user interfaces? The scope encompasses both theoretical foundations of trust in professional human-AI interactions and practical technical implementations that enterprise software developers can deploy. By examining these questions, this article aims to provide a comprehensive framework for trust-building that accommodates the complex requirements of enterprise settings.

## 2. Theoretical Framework: Understanding Trust in Human-AI Interactions

Trust forms the essential foundation for effective human-AI collaboration in enterprise environments. As organizations increasingly integrate AI assistants into their workflows, establishing a comprehensive understanding of how trust manifests in these interactions becomes crucial. This section examines the theoretical underpinnings of trust in enterprise AI contexts, exploring its definition, components, relationship to responsible AI principles, and connections to existing research on human trust in automated systems.

### 2.1. Defining Trust in the Context of Enterprise AI

Trust in enterprise AI extends beyond simple confidence in system accuracy to encompass a multidimensional relationship between users and AI assistants. In enterprise contexts, trust can be defined as the willingness of organizational users to rely on AI systems for consequential tasks, to accept AI-generated information in decision-making processes, and to share sensitive information with these systems. This willingness emerges from users' beliefs about the system's capabilities, intentions, and governance. The IEEE AIS Trust and Agency Committee emphasizes that trust in AI differs fundamentally from interpersonal trust, noting that it involves "the establishment of reliance on AI systems with an understanding of their capabilities and limitations" [4]. In enterprise settings, this trust operates at both individual and organizational levels, with institutional trust mechanisms often conditioning individual users' trust perceptions. This multifaceted nature makes trust particularly complex in enterprise environments where professional responsibilities, organizational policies, and individual preferences intersect.

### 2.2. Components of AI Trust: Capability, Reliability, Transparency, and Data Security

Trust in enterprise AI assistants comprises several interconnected components that collectively shape user perceptions and behaviors. First, capability trust refers to users' confidence in the AI system's fundamental competence to perform its intended functions across relevant domains. Second, reliability trust encompasses consistency in performance, predictability of behavior, and availability when needed. Third, transparency trust relates to the system's explainability, including clear communication about its capabilities, limitations, and the sources of its information. Fourth, data security trust involves confidence that the AI system will handle sensitive information appropriately, respect privacy boundaries, and maintain information confidentiality. Chen et al. note that these components interact dynamically, with deficiencies in one area potentially undermining trust across all dimensions [3]. Additionally, these components

manifest differently across various enterprise contexts, with regulated industries often placing heightened emphasis on transparency and data security aspects of trust.

**Table 1** Core Components of Enterprise AI Trust and Their Technical Implementations [3, 4]

| Trust Component | Description | Key Technical Implementation Approaches |
|---|---|---|
| Capability Trust | User confidence in AI system's competence | Performance monitoring, confidence scoring |
| Reliability Trust | Consistency in performance and predictability | Availability metrics, failure handling mechanisms |
| Transparency Trust | System explainability and clear communication | AI identity labeling, source attribution |
| Data Security Trust | Appropriate handling of sensitive information | Encryption, access controls, data minimization |

## 2.3. The Role of Responsible AI Principles in Trust-Building

Responsible AI principles provide an ethical foundation for developing trustworthy enterprise AI systems. These principles—typically including fairness, accountability, transparency, explainability, and privacy—establish normative standards that guide both technical implementation and organizational governance. Chen et al. argue that responsible AI principles must transition "from principles to practice" through concrete implementation mechanisms that users can perceive and experience [3]. In enterprise contexts, responsible AI principles often extend to include additional considerations such as business continuity, regulatory compliance, and alignment with organizational values. These principles serve multiple trust-building functions: they provide developers with ethical guardrails during system creation, offer users assurance that systems were designed with appropriate considerations, and establish shared expectations between organizations and their stakeholders regarding AI system behavior. When properly implemented, responsible AI principles create a virtuous cycle where ethical design increases trust, which in turn encourages appropriate system use and further ethical refinement.

## 2.4. Review of Existing Literature on Human Trust in Automated Systems

Research on human trust in automated systems provides valuable insights that inform understanding of enterprise AI trust dynamics. This literature spans multiple disciplines including human-computer interaction, organizational psychology, information systems, and cognitive science. Early work focused primarily on automation in high-risk environments such as aviation and industrial control systems, establishing fundamental principles regarding overtrust, undertrust, and appropriate reliance. More recent research has examined trust in advisory systems, algorithm aversion in decision-making contexts, and the impact of anthropomorphism on trust perceptions. The IEEE AIS Trust and Agency Committee highlights the importance of "enabling end-user agency" alongside trust, noting that genuine trust requires users to maintain meaningful control and understanding of AI systems [4]. Enterprise AI trust research builds upon these foundations while addressing unique organizational factors including professional identity, institutional policies, and complex stakeholder relationships. This growing body of literature increasingly recognizes trust not as a static property but as a dynamic relationship that evolves through ongoing interactions between users, AI systems, and organizational contexts.

# 3. Identity and Information Verification Mechanisms

As enterprise AI assistants assume increasingly central roles in workplace software, mechanisms for clearly identifying AI-generated content and verifying information sources become essential components of trust architecture. This section examines technical approaches to signaling AI identity, tracking information provenance, and establishing verification systems that support appropriate user reliance on AI-provided information.

## 3.1. AI Labeling Standards and Implementation Approaches

AI labeling systems provide explicit indicators that content or interactions originate from artificial intelligence rather than human sources. These signals serve multiple trust functions: they establish appropriate user expectations, prevent deception, and create transparency about the nature of the interaction. Wittenberg et al. note that effective AI labeling systems must balance visibility with user experience considerations, as overly intrusive labels may disrupt workflow while subtle indicators might go unnoticed [5]. Technical implementation approaches vary considerably, ranging from

persistent visual markers and distinctive interaction patterns to metadata tagging and digital watermarks. In enterprise contexts, these labeling systems often need to function across diverse software environments including chat interfaces, document systems, collaboration tools, and workflow applications. Standardization efforts have begun emerging to create consistent labeling approaches across platforms, though significant fragmentation remains in implementation practices. The technical architecture for AI labeling typically requires coordination across multiple system layers, including model outputs, middleware processing, and front-end rendering components.

**Table 2** Enterprise AI Labeling Implementation Approaches [5, 6]

| Labeling Approach | Implementation Mechanism | Enterprise Application |
|---|---|---|
| Visual Indicators | Distinctive UI elements | Document management systems |
| Interaction Patterns | Distinctive conversation styles | Customer service platforms |
| Metadata Tagging | Machine-readable indicators | Content management systems |
| Digital Watermarking | Embedded signals in content | Enterprise communication platforms |

## 3.2. Citation Systems: Technical Architecture for Source Tracking

Citation systems provide technical infrastructure for tracking and communicating the sources of information provided by AI assistants. These systems enable users to verify claims, assess information credibility, and understand the contextual basis for AI-generated content. The technical architecture for enterprise citation systems typically comprises several interconnected components: content indexing systems that maintain relationships between information fragments and their sources, retrieval mechanisms that associate AI outputs with relevant source materials, citation generation processes that format reference information appropriately, and front-end components that render citations in user interfaces. Arooj emphasizes that effective source tracking in enterprise environments requires systems capable of maintaining provenance across multiple information transformations, including paraphrasing, summarization, and synthesis of multiple sources [6]. Implementation challenges include tracking information lineage through complex processing pipelines, handling proprietary or access-restricted source materials, and addressing the computational overhead of maintaining detailed provenance records. Additionally, enterprise citation systems must often integrate with existing organizational knowledge management infrastructure while supporting industry-specific citation formats and standards.

## 3.3. Information Provenance Verification Methods

Information provenance verification extends beyond simple citation to establish confidence in the authenticity, accuracy, and appropriateness of AI-provided information. These verification mechanisms serve as technical safeguards against misinformation, hallucination, and inappropriate information disclosure. Enterprise implementations typically employ layered verification approaches that combine multiple methods. Source validation techniques confirm the legitimacy and reliability of information origins, often leveraging organizational content authority systems or external validation services. Content verification mechanisms assess the accuracy of specific claims through fact-checking processes, comparison with trusted sources, or consistency analysis. Context verification evaluates whether information is appropriate for the specific organizational setting, user role, and task context. Wittenberg et al. observe that verification systems must balance thoroughness with performance considerations, as extensive verification processes may introduce latency that diminishes user experience [5]. Advanced implementations increasingly employ verification strategies that adjust dynamically based on information sensitivity, task criticality, and confidence levels, applying more rigorous verification to high-stakes contexts.

## 3.4. Case Studies of Successful Identity Signaling in Enterprise Systems

Implementations of identity and verification mechanisms across various enterprise contexts provide valuable insights into effective technical approaches. Financial service organizations have developed AI assistant systems featuring persistent visual indicators that maintain consistent identity signals across multiple interaction modalities, integrated with citation systems that provide contextual links to regulatory documentation and organizational policies. Healthcare enterprises have implemented verification systems that combine real-time fact-checking with clear sourcing of medical information, using standardized citation formats familiar to clinical professionals. Professional services firms have deployed AI assistants that employ distinctive conversation patterns explicitly communicating AI identity while providing interactive access to source materials supporting recommendations. Across these implementations, Arooj notes that successful approaches share several characteristics: they integrate identity and verification mechanisms throughout the entire user experience rather than treating them as afterthoughts, they align technical implementation

with existing organizational trust signals, and they maintain consistent identity across different interface contexts [6]. These case studies demonstrate that effective identity and verification systems must be tailored to specific enterprise contexts while adhering to emerging standards and best practices.

## 4. Data Handling Safeguards for Enterprise Contexts

Enterprise AI assistants routinely process, analyze, and generate content based on organizational data that varies significantly in sensitivity and regulatory requirements. Establishing appropriate safeguards for this data handling represents a critical dimension of trust architecture. This section examines technical approaches to integrating AI systems with enterprise data classification frameworks, implementing sensitivity detection and handling protocols, addressing compliance requirements in regulated industries, and preserving privacy in AI interactions.

### 4.1. Integration with Enterprise Data Classification Frameworks

Enterprise data classification frameworks provide structured approaches to categorizing information based on sensitivity, regulatory requirements, and organizational value. For AI assistants to function as trusted participants in enterprise workflows, they must align with these existing classification structures. Technical integration typically occurs at multiple levels: AI systems must recognize and interpret classification labels from content management systems, respect access control boundaries based on classification levels, and maintain appropriate classification metadata throughout processing pipelines. Gluckd and Robmazz emphasize that effective integration requires AI systems to "understand the downstream impacts of classification decisions" rather than simply recognizing labels [7]. Implementation approaches vary based on organizational maturity, ranging from basic recognition of explicit classification tags to sophisticated systems that maintain classification lineage through complex transformations. Technical challenges include handling inconsistent or missing classification metadata, addressing classification conflicts across federated systems, and maintaining classification integrity during AI processing. Advanced implementations increasingly employ machine learning techniques to predict likely classification levels for unlabeled content based on semantic similarity to previously classified materials.

### 4.2. Sensitivity Detection and Handling Protocols

While classification frameworks provide structural guidance for data handling, additional mechanisms for sensitivity detection and handling are necessary to address content that may be unclassified, misclassified, or created during AI interactions. These systems serve as technical safeguards against inappropriate disclosure or processing of sensitive information. Pattern-based detection identifies known sensitive data types such as personal identifiers, financial information, and protected health information through regular expressions, entity recognition, and heuristic rules. Semantic sensitivity detection employs machine learning approaches to identify potentially sensitive content based on contextual understanding, even when explicit patterns are absent. Response protocols determine appropriate system actions when sensitive content is detected, ranging from applying classification labels and limiting distribution to implementing special handling procedures or declining certain processing requests. Feretzakis et al. note that sensitivity detection systems face significant challenges including cultural variations in sensitivity perceptions, domain-specific sensitivity types, and evolving definitions of sensitive information [8]. Enterprise implementations increasingly employ layered approaches combining multiple detection methods with graduated response protocols that balance security requirements with usability considerations.

### 4.3. Compliance Considerations for Regulated Industries

Organizations in regulated industries face additional requirements for AI data handling related to specific legal and regulatory frameworks. These requirements introduce distinct technical challenges for establishing trusted AI interactions. Healthcare organizations must implement AI systems that maintain HIPAA compliance through appropriate de-identification techniques, access controls, and audit trails. Financial services firms require AI assistants that adhere to requirements including proper handling of non-public information, maintenance of appropriate information barriers, and capabilities for supervisory review. Public sector organizations need systems that support Freedom of Information Act requirements, records management regulations, and appropriate handling of controlled unclassified information. Technical implementations typically involve multiple compliance mechanisms including data residency controls that maintain appropriate geographic boundaries for processing and storage, retention management capabilities that support regulatory timeframes, and compliance metadata that tracks relevant regulatory frameworks for specific content. Gluckd and Robmazz emphasize that compliance controls should be "integrated by design rather than applied as overlays" to ensure effectiveness within production systems [7]. As regulatory requirements continue evolving to address AI specifically, technical implementations must maintain flexibility to incorporate emerging compliance obligations.

## 4.4. Technical Approaches to Privacy Preservation in AI Interactions

Privacy preservation in AI interactions extends beyond basic security measures to include sophisticated technical approaches that protect sensitive information while maintaining system utility. These techniques aim to establish trust by demonstrating that AI assistants can provide value without unnecessary exposure of private data. Federated learning approaches enable model improvement without centralizing sensitive data by training models across distributed datasets and sharing only model updates rather than raw information. Differential privacy introduces controlled statistical noise to prevent extraction of individual data points while preserving aggregate insights. Secure multi-party computation allows multiple entities to jointly compute functions over inputs while keeping those inputs private. Homomorphic encryption permits computation on encrypted data without decryption. Feretzakis et al. observe that enterprise implementations often combine multiple privacy-preserving techniques in layered architectures tailored to specific use cases and sensitivity levels [8]. Implementation challenges include performance impacts of privacy-preserving techniques, complexity of properly configuring privacy parameters, and difficulties in communicating privacy properties to users in understandable terms. Advanced implementations increasingly employ context-aware privacy mechanisms that adjust protection levels based on interaction context, data sensitivity, and user authorization.

**Table 3** Privacy-Preserving Techniques for Enterprise AI Systems [7, 8]

| Technique | Working Principle | Enterprise Applications |
|---|---|---|
| Federated Learning | Training across distributed datasets | Multi-regional enterprises |
| Differential Privacy | Adding controlled statistical noise | Financial analytics |
| Secure Multi-party Computation | Joint computation with private inputs | Cross-organizational collaboration |
| Homomorphic Encryption | Computation on encrypted data | Healthcare AI |

## 5. System Architecture for Trusted AI Workflows

As enterprise AI assistants become more sophisticated, they increasingly operate as components within larger systems rather than standalone applications. These multi-agent workflows, where tasks transfer between different AI subsystems or human collaborators, introduce unique trust challenges related to context preservation, authorization boundaries, and interface consistency. This section examines architectural approaches for maintaining trust throughout these complex workflows, exploring secure context transfer mechanisms, authentication and authorization models, API design principles, and implementation requirements across the technology stack.

### 5.1. Secure Context Transfer in Multi-Agent Systems

Multi-agent AI workflows require secure and comprehensive transfer of context between system components to maintain coherent user experiences and appropriate handling of information. Context transfer encompasses multiple dimensions including conversation history, user intent, task status, data references, and security parameters. Wu and He emphasize that secure context transfer systems must maintain resilience against various threat vectors including information tampering, context corruption, and timing attacks that exploit transfer delays [9]. Technical implementation approaches vary based on architectural complexity, ranging from basic serialization with encryption to sophisticated context management services that maintain state across distributed systems. Key technical challenges include minimizing latency during context transfers, preserving semantic integrity across different agent capabilities, and maintaining appropriate security boundaries as context moves between systems with different trust characteristics. Enterprise implementations increasingly employ intermediate context representation formats that standardize information across heterogeneous systems while incorporating cryptographic mechanisms to verify context integrity and provenance throughout transfer processes.

### 5.2. Authentication and Authorization Models for AI Handoffs

Authentication and authorization mechanisms ensure that AI assistants maintain appropriate permissions and identity verification throughout complex workflows involving multiple systems. These mechanisms establish clear boundaries of authority and access during agent transitions. South et al. describe a framework for "authenticated delegation and authorized AI agents" that maintains security properties throughout task handoffs [10]. Authentication approaches verify the identity and integrity of AI systems receiving delegated tasks, often employing techniques such as system certificates, secure tokens, and cryptographic signatures. Authorization frameworks define and enforce permissions boundaries when tasks transfer between agents, determining what data each system can access and what actions it can

perform. Enterprise implementations typically incorporate concepts including least privilege principles that limit each agent to minimum necessary permissions, role-based access controls that align AI permissions with organizational structures, and dynamic authorization adjustments based on context and task requirements. Technical challenges include maintaining authorization continuity across system boundaries, communicating permission constraints in machine-readable formats, and implementing appropriate fallback mechanisms when authorization requirements cannot be satisfied.

## 5.3. API Design Principles for Maintaining Trust Across Transitions

Application Programming Interfaces (APIs) serve as the connective tissue in multi-agent AI workflows, defining how different systems exchange information and functionality. API design significantly influences trust preservation during transitions between AI subsystems. South et al. identify several key API design principles for trusted transitions including explicitness in representing trust parameters, consistency in error handling across boundaries, and verifiability of cross-system information flows [10]. Technical implementation patterns include trust contracts that formally specify expectations between systems, standardized trust metadata that communicates security and provenance information across interfaces, and trust verification mechanisms that validate adherence to specified parameters. Additional considerations include designing appropriate rate limiting and throttling to prevent abuse, implementing comprehensive logging and auditing capabilities, and establishing clear versioning approaches that maintain compatibility while enabling security improvements. Enterprise implementations increasingly employ API governance frameworks that ensure consistent application of trust principles across organizational interfaces, often incorporating automated compliance checks and security scanning as part of API lifecycle management.

## 5.4. Implementation Requirements Across the Technology Stack

Implementing trusted multi-agent workflows requires coordinated approaches across multiple layers of the technology stack, from infrastructure components to user experience elements. At the infrastructure layer, secure communication channels establish encrypted, authenticated connections between system components, while containerization and isolation techniques maintain appropriate boundaries between processing environments. Database and storage systems implement fine-grained access controls and maintain tamper-evident logs of context transfers. Middleware components provide identity management, permissions enforcement, and context synchronization across distributed systems. Application services implement domain-specific trust preservation logic, including sensitivity detection, classification handling, and compliance validation. User interface components communicate workflow transitions transparently while maintaining consistent trust signals throughout handoff processes. Wu and He note that successful implementations require "coordinated design across all system layers" rather than addressing trust considerations in isolation [9]. Implementation challenges include maintaining consistent trust properties across systems with different technical architectures, balancing security requirements with performance considerations, and coordinating trust implementations across organizational boundaries in enterprise settings that involve multiple vendors and internal teams.

# 6. Conclusion

The technical mechanisms for building user trust in enterprise AI assistants represent essential infrastructure for the successful integration of artificial intelligence into professional environments. AI labeling, citation systems, sensitivity handling protocols, and secure handoff architectures collectively create a foundation for appropriate reliance on these increasingly sophisticated tools. This interconnected trust framework addresses the multidimensional requirements of enterprise contexts, balancing transparency with efficiency, security with accessibility, and standardization with contextual adaptation. Organizations implementing these technical safeguards can expect increased user confidence, more widespread adoption, and reduced resistance to AI integration within business-critical workflows. While technical approaches alone cannot establish complete trust—organizational culture, governance structures, and individual experiences remain crucial factors—they provide the necessary conditions for trustworthy interactions to develop and persist. As enterprise AI assistants continue evolving from narrow task-specific tools toward general-purpose collaborative systems, these trust mechanisms will require ongoing refinement to accommodate new capabilities, emerging threats, and changing regulatory landscapes. The future direction of this field points toward increasingly standardized trust architectures that function across organizational boundaries while maintaining the flexibility to address industry-specific requirements and unique organizational contexts.

## References

[1]     Cleidson Ronald Botelho de Souza, et al., "The Fine Balance Between Helping With Your Job and Taking It," IEEE Software, Nov/Dec 2024. https://www.cs.ubc.ca/~bestchai/papers/ieee-software24-ai-code-assistants.pdf

[2]     Christopher Bull, Ahmed Kharrufa, "Generative Artificial Intelligence Assistants in Software Development Education: A Vision for Integrating Generative Artificial Intelligence Into Educational Practice, Not Instinctively Defending Against It," IEEE Software, 08 August 2023. https://ieeexplore.ieee.org/document/10213396/authors#authors

[3]     Fang Chen, et al., "Artificial Intelligence Ethics and Trust: From Principles to Practice," IEEE Intelligent Systems, 13 November 2023. https://ieeexplore.ieee.org/stampPDF/getPDF.jsp?arnumber=10352170

[4]     IEEE AIS Trust and Agency Committee, "Enabling End-User Agency and Trust in Artificial Intelligence Systems," IEEE Standards Association, 26 October 2021. https://standards.ieee.org/beyond-standards/industry/technology-industry/enabling-end-user-trust-in-artificial-intelligence-in-the-algorithmic-age/

[5]     Chloe Wittenberg, et al., "Labeling AI-Generated Content: Promises, Perils, and Future Directions," MIT Schwarzman College of Computing, 28 November 2023. https://computing.mit.edu/wp-content/uploads/2023/11/AI-Policy_Labeling.pdf

[6]     Arooj, "Effective Source Tracking in Retrieval-Augmented Generation (RAG) Systems," Chitika AI Research, 28 March 2025. https://www.chitika.com/source-tracking-rag/

[7]     Gluckd and Robmazz, "Create a Well-Designed Data Classification Framework," Microsoft Learn, 20 June 2024. https://learn.microsoft.com/en-us/compliance/assurance/assurance-create-data-classification-framework

[8]     Georgios Feretzakis, et al., "Privacy-Preserving Techniques in Generative AI and Large Language Models," MDPI, 4 November 2024. https://www.mdpi.com/2078-2489/15/11/697

[9]     Yiming Wu, Xiongxiong He, "Secure Consensus Control for Multi-Agent Systems With Attacks and Communication Delays," IEEE/CAA Journal of Automatica Sinica, January 2017. https://ieeexplore.ieee.org/stampPDF/getPDF.jsp?arnumber=7815560

[10]    Tobin South, et al., "Authenticated Delegation and Authorized AI Agents," arXiv (Computers and Society), 16 January 2025. https://arxiv.org/abs/2501.09674