



# Comparative Analysis: Enterprise AI Agent Governance vs Zero Trust Agent Frameworks

## Introduction

Autonomous AI agents are emerging as powerful tools in enterprise workflows, but they also introduce new governance and security challenges <sup>1</sup>. Two internal framework concepts address these needs from different angles:

- **Enterprise AI Agent Governance Framework (EAAGF):** A structured model for deploying, monitoring, and controlling AI agents in regulated enterprises, including a maturity model, layered controls (auditability, supervision, autonomy), and risk mitigation scoring.
- **Zero Trust Authentication Framework for Autonomous Agents:** A taxonomy and methodology for ensuring identity, fine-grained access control, and continuous behavioral trust evaluation for AI agents, especially as they operate across organizational boundaries.

This report compares the novelty and competitive landscape of these two frameworks, examining existing similar frameworks (academic, industry, open-source), their uniqueness, related efforts, and potential venues and communities for publication or adoption. A summary table highlights each framework's novelty and overlap with prior work.

## Enterprise AI Agent Governance Framework (EAAGF)

**Description:** EAAGF is envisioned as a comprehensive governance model for AI agents in enterprise settings. It would define **maturity levels** of agent deployment (from basic human-supervised agents up to highly autonomous agents) and enforce **control layers** – for example, ensuring **auditability** of agent decisions, human **supervision** or oversight triggers at appropriate points, and bounded **autonomy** levels. The framework also proposes a **risk mitigation scoring** system to quantify the operational risks of an agent and ensure that higher-risk agent actions face stricter controls or approvals.

**Existing Similar Frameworks:** Several emerging frameworks echo parts of EAAGF's scope, though none appear to combine all its elements in one model. Notable related efforts include:

- **Agentic AI Governance Concepts:** Cybersecurity and AI firms have highlighted the need for governance beyond traditional AI model oversight. For instance, Palo Alto Networks defines *agentic AI governance* as “*the structured management of delegated authority in autonomous AI systems*” with explicit boundaries on what agents can do and who is accountable <sup>1</sup>. This underscores the shift from governing static model outputs to governing **AI actions and decisions** in runtime <sup>2</sup>. EAAGF aligns with this philosophy by enforcing *operational control* over agent actions (not just policy compliance at design time) <sup>3</sup>.

- **NIST AI RMF and ISO 42001:** Broad AI governance frameworks exist, such as NIST's AI Risk Management Framework and the new ISO/IEC 42001 standard for AI management systems. These provide high-level guidance on AI risk identification, assessment, and lifecycle management <sup>4</sup> <sup>5</sup>. However, they are not tailored to autonomous agents' unique challenges – for example, NIST AI RMF emphasizes mapping and mitigating risks in AI systems generally, but “*most AI governance programs were designed for model outputs, not autonomous actions*” <sup>2</sup>. EAAGF would extend such frameworks with agent-specific controls (like real-time oversight and action constraints) that current standards only imply.
- **Autonomy Levels and Maturity Models:** The concept of graded autonomy or maturity is emerging. The Cloud Security Alliance (CSA) has discussed a **six-level autonomy taxonomy** for agentic AI, inspired by self-driving car levels <sup>6</sup>. Others propose simpler tiered models (e.g. *human-in-control, bounded autonomy, full autonomy*) to classify how much independence an agent has <sup>6</sup>. Similarly, some industry frameworks include maturity models for agent governance. For example, Accuvity's *Agent Integrity Framework* defines five maturity levels across capabilities – from Level 1 (legacy controls only) up to Level 5 (full real-time enforcement of policies and intent) <sup>7</sup>. EAAGF's maturity model would be comparable, guiding enterprises from basic controls to advanced, integrated governance for agents.
- **Multi-Agent Oversight Frameworks:** With enterprises deploying *multi-agent systems*, experts note that governance must cover agent interactions and emergent behaviors. One source suggests organizations must “*extend governance frameworks to address agent-to-agent communication protocols, coordination mechanisms, and collective decision-making... including clear autonomy boundaries and human oversight triggers when agents collaborate on high-stakes decisions.*” <sup>8</sup>. EAAGF directly addresses this by incorporating supervision layers and defined autonomy limits as part of its control stack. This focus on multi-agent oversight differentiates it from traditional AI governance that usually considers single models in isolation.

**Novelty and Differentiation:** EAAGF appears **partially novel** – it builds upon known concepts (AI governance, risk management, autonomy levels) but integrates them into a specialized framework for enterprise agents. The idea of a **maturity model** for agent governance is relatively new; while maturity matrices exist for general AI adoption, few are tailored to autonomous agent controls. The inclusion of specific **control layers** (audit logs, supervision checkpoints, adjustable autonomy) is a differentiator, ensuring not just policy compliance but live agent *command and control*. Additionally, a quantitative **risk scoring** mechanism for agent behaviors would be an innovative contribution – current frameworks like NIST AI RMF are qualitative or process-oriented, whereas EAAGF's scoring could provide ongoing risk metrics per agent or task (comparable to how financial firms score model risk). This could help regulated industries dynamically enforce stricter governance on higher-risk AI agents. In summary, similar ideas are in the air, but EAAGF's **comprehensive, structured approach** (combining governance processes with technical control layers and metrics) would be distinct in the market.

**Closest Comparable Efforts:** Key organizations and efforts related to EAAGF include:

- **Cloud Security Alliance (CSA):** Through its AI Safety Initiative, CSA is exploring governance and safety for agentic AI (e.g., autonomy taxonomies and AI control frameworks). Though CSA's focus often leans toward security, their work on **AI Controls** and a potential autonomy-level framework

suggests alignment with EAAGF's aims. Collaboration or cross-referencing with CSA's upcoming guidelines (like the CSA AICM – AI Controls Matrix) could validate EAAGF's structure.

- **Consulting & Research Firms:** McKinsey, Gartner, and others have begun advising enterprises on agent governance. McKinsey's playbook for agentic AI security notes that enterprises need "*standardized oversight processes, defined accountability for agent actions, and triggers for escalation*" as part of governance <sup>9</sup>. It also highlights that only 1% of organizations felt their AI (let alone agent) adoption was mature <sup>10</sup>, pointing to a gap EAAGF can fill.
- **Enterprise AI Platforms:** Companies like IBM, Microsoft, and Salesforce (e.g., Salesforce's "Agentforce" concept) are embedding governance into their AI agent offerings. For example, IBM's AI Governance toolkit or Microsoft's Azure OpenAI Service include features for **auditability** and **compliance** – EAAGF could position as a vendor-neutral framework that complements such platform-specific controls with an overarching maturity model.
- **Open-Source Initiatives:** There is currently no dominant open-source "agent governance framework," but projects in the AI safety community (like toolkits for **AI alignment** or **human-in-the-loop orchestration**) partially overlap. EAAGF could be the first to codify best practices from these communities (e.g., reinforcement learning with human feedback, audit logging libraries, etc.) into a unified framework.

**Positioning Strategies:** Given some overlap with existing ideas, EAAGF should be positioned as an **integrated and enterprise-ready** framework. Emphasize that it goes beyond high-level principles by providing a **practical maturity roadmap** and **layered controls** specifically for AI agents. For instance, unlike general AI governance guidelines, EAAGF would tell a bank or hospital *exactly* how to restrict an agent's autonomy based on risk scoring, and how to evolve those controls as the organization's capabilities mature. If similar frameworks exist (e.g., Acuity's security-focused model), EAAGF can be differentiated by its **holistic scope** – covering not just security integrity but also ethical compliance, performance monitoring, and business value alignment in one maturity model. In essence, EAAGF can be the "**COBIT for AI agents**" or "**CMMI for autonomous AI**", filling a gap between broad AI principles and low-level system configs.

## Zero Trust Authentication Framework for Autonomous Agents

**Description:** The Zero Trust Authentication Framework for Autonomous Agents aims to apply **Zero Trust Architecture (ZTA)** principles to the world of AI agents. In a zero-trust model, no agent (even an internal one) is inherently trusted; every interaction must be authenticated, authorized, and continuously validated. This framework would establish a **taxonomy** for agent identity and trust, and a methodology to enforce:

- **Robust Agent Identity:** Defining a unique, verifiable identity for each AI agent (potentially using technologies like decentralized identifiers). Each agent is treated as a "non-human user" with its own credentials and attributes, rather than an invisible part of a software process <sup>11</sup> <sup>12</sup>. This includes managing the agent's identity lifecycle (issuance, rotation, revocation) just as one would for human accounts <sup>13</sup>.
- **Fine-Grained Access Control:** Implementing strict **least privilege** and context-based access for agents. The framework would use mechanisms like role-based or attribute-based access control (RBAC/ABAC) and just-in-time credentials so an agent can only call APIs or resources authorized for its current task <sup>14</sup> <sup>15</sup>. For example, an agent gets a short-lived token to read one database record if needed, rather than broad database read access.

- **Behavioral Trust Evaluation:** Continuously monitoring agent behavior and environment to adjust trust levels. The framework might include a **trust scoring** engine that raises or lowers an agent's trust based on its actions, history, and compliance with policies <sup>16</sup>. Suspicious or out-of-policy behavior (e.g., an agent accessing unusual data or running unexpected tool commands) would trigger containment or additional authentication challenges in real time.
- **Cross-Boundary Verification:** Ensuring that when agents operate across enterprise boundaries (e.g. an agent from Company A requesting data from Company B), there is a federated trust mechanism. This could involve **mutual authentication** protocols, verifiable claims about the agent's identity and policy compliance, and audit trails shared between parties. The goal is to prevent "rogue" or impersonated agents from exploiting inter-company interactions – essentially extending zero trust network concepts to agent-to-agent API ecosystems.

**Existing Similar Frameworks:** The notion of applying zero trust to AI agents is cutting-edge, but several parallel efforts are exploring it:

- **CSA's Zero Trust for Agentic AI:** The Cloud Security Alliance published a paper and blog proposing a "Zero-Trust IAM framework, designed specifically for agents in distributed ecosystems." It outlines core principles like "Never Trust, Always Verify" every agent action, using **DIDs (Decentralized IDs)** and **Verifiable Credentials (VCs)** for dynamic agent identities, and **continuous monitoring** of agent behavior <sup>17</sup> <sup>18</sup>. Key components of their approach include a cryptographically anchored agent identity, an Agent Name Service for discovery, just-in-time access tokens, and continuous **trust scoring based on behavior and compliance** <sup>16</sup>. This closely mirrors the intent of the Zero Trust Authentication Framework. In fact, an academic paper (Huang *et al.* 2025) on arXiv details a similar architecture: a "novel Agentic AI IAM framework" using DIDs/VCs, fine-grained access control, unified session management, and zero-knowledge proofs for policy compliance <sup>19</sup>. These sources indicate that the concept is active in the research community, with CSA and collaborators pushing for standards.
- **Enterprise Identity Solutions (Okta, HashiCorp, etc.):** Identity and security companies have begun adapting zero trust to non-human identities. Okta, for example, advocates treating "every agent as a first-class identity" and enforcing **Zero Trust (never assume one agent trusts another)**, with micro-segmentation of agent permissions <sup>20</sup>. They emphasize practices like short-lived credentials, per-request verification, and isolating agent credentials to prevent lateral movement <sup>14</sup> <sup>21</sup>. HashiCorp has similarly discussed managing machine and agent identities at scale with zero trust principles (dynamic secrets, identity-based authentication) <sup>22</sup>. The Zero Trust Agent Framework would formalize these best practices into a cohesive taxonomy and reference model.
- **Agent Security Frameworks:** Beyond identity, some frameworks focus on agent behavior integrity (which complements zero trust). For instance, Acuity's Agent Integrity Framework (mentioned earlier) introduces the idea of "*continuous verification at the semantic level*" – confirming an agent *should* take a given action, not just that it *can* <sup>23</sup> <sup>24</sup>. One pillar is **Identity & Attribution**, ensuring every agent-initiated action is traceable to a verified agent identity and user request <sup>25</sup>. Another is **Behavioral Consistency**, detecting when an agent deviates from its expected usage patterns <sup>26</sup>. These align with zero trust's call for constant evaluation. The Zero Trust Authentication Framework could incorporate such concepts (intent verification, anomaly detection) into its trust evaluations, going beyond traditional access control.

- **Inter-Agent Communication Protocols:** Big tech companies are recognizing the need for secure agent interoperability. Efforts like **Anthropic's "Model Context Protocol," Cisco's "Agent Connect," Google's "Agent2Agent,"** and IBM's agent communication standards are “*under development but not yet mature*” <sup>27</sup>. All of them grapple with authenticating and permissioning agents in multi-agent workflows. The proposed framework could leverage or inform these protocols, ensuring a taxonomy where any agent communication involves mutual auth, logging, and policy enforcement at the boundaries. In essence, it would provide the zero trust rulebook that such protocols should follow (e.g., requiring token exchange, verifying agent identity claims, etc.).

**Novelty and Differentiation:** The Zero Trust Authentication Framework for Autonomous Agents is **highly novel**, though it converges with ideas only now appearing in late 2025 and 2026. Zero trust principles themselves are well-known in cybersecurity, but applying them to AI agent ecosystems is nascent. A few differentiators and gaps it addresses:

- **Non-Human IAM Paradigm:** Traditional IAM struggles with the “*dynamic, interdependent, and ephemeral nature*” of AI agents <sup>28</sup>. Our framework explicitly tackles this by redefining identity (e.g. providing an Agent ID construct with rich attributes about the agent’s capabilities and context). This goes beyond the coarse identities (API keys or service accounts) used today. While the CSA paper and others have started defining *Agent IDs*, a comprehensive taxonomy that can be adopted across industries would be novel.
- **Unified Trust Scoring:** Continuous trust or risk scoring for identities is an emerging concept. Some zero trust implementations use session risk scores for humans; for agents, the framework proposes scoring based on policy compliance, anomaly detection, and even external factors. CSA’s architecture mentions a “*trust engine*” and dynamic adjustment of privileges <sup>29</sup> <sup>30</sup>, which is on point. However, the novelty could lie in formalizing **what factors contribute to trust** for an AI agent (e.g., model confidence, alignment checks, past incidents) and how to quantitatively evaluate them. Few existing solutions do this in a standardized way.
- **Cross-Enterprise and Decentralized Trust:** Many current discussions (Okta, etc.) assume an enterprise’s internal environment. The proposed framework emphasizes agents operating *across enterprise boundaries*, which likely requires a **decentralized or federated trust model** (akin to how SAML or OAuth federates human identity). Introducing decentralized identifiers and verifiable credentials for agents would be cutting-edge – effectively giving agents “digital passports” that any organization can verify. The arXiv paper by Huang *et al.* hints at this via an Agent Naming Service and global policy layer <sup>16</sup>. Our framework stands to be among the first that could be implemented in open-source (e.g., using blockchain or PKI for agent credentials) to allow inter-company agent trust without a single authority. This cross-org focus differentiates it from most vendor-specific zero trust solutions.

In summary, the framework’s novelty is not in the zero trust concept itself, but in **extending it to AI agents with a full identity & trust lifecycle**. It unifies ideas from identity management, AI behavior monitoring, and secure multi-agent protocols into one methodology. While parallel projects (CSA, academic research) exist, this framework could be positioned as an early **reference architecture or open standard** that others can implement, much like the Zero Trust Architecture model (NIST SP 800-207) did for general IT.

**Closest Comparable Efforts:** Key players and efforts in this space include:

- **Cloud Security Alliance (CSA):** As noted, CSA is at the forefront (their Zero Trust Advancement Center, AI Safety Working Group, and papers on agent identity and MAESTRO threat model). Collaborating with CSA or aligning terminology (DID, VC, trust fabric layers) could lend credibility.
- **OWASP & Security Community:** The OWASP Foundation's *GenAI Security* project published an "*OWASP Top 10 for Agentic AI Security*" <sup>31</sup>, indicating security experts are already rallying around agent-specific threats. They and conference communities (DEF CON AI Village, Black Hat briefings) are primed for zero trust solutions for AI – indeed, a *Black Hat 2025* demo showed prompt injection attacks that traditional controls missed <sup>32</sup>. Our framework will resonate as a defense strategy in such forums.
- **Tech Companies & Consortia:** Identity-focused companies like **Okta**, **Microsoft Entra**, **HashiCorp**, certificate authorities like **Keyfactor** (which is exploring machine identity for AI <sup>33</sup>), and startups like **Xage** (zero trust for IoT/AI) are all potential collaborators or competitors. IBM's and Google's nascent protocols for agent comms and Anthropic's safety mechanisms mean those firms have interest in secure agent interoperability. The framework could position itself as a unifying layer that these protocols plug into – potentially gaining traction if presented in standards bodies or open-source projects that these companies contribute to.

**Positioning Strategies:** To stand out amid these efforts, this Zero Trust Agent Framework should highlight its **comprehensiveness and practical blueprint**. For example, it can be framed as "*the Zero Trust playbook for AI Agents*," combining identity management, access policies, and trust analytics. If parts of the idea overlap with CSA's proposals, a strategy could be to publish collaboratively (or at least reference their work) to gain community buy-in rather than appear redundant. Emphasize unique aspects like cross-enterprise trust and **behavioral authentication** (continuously verifying an agent by what it's doing, not just who it is). Also, positioning it as an **open framework** (with perhaps a GitHub reference implementation or schema) could differentiate it from proprietary solutions and encourage adoption in developer communities building multi-agent systems. Given rising concern about "rogue AI" in enterprises, a clear, implementable framework that "*never assumes agent trust by default*" and provides measurable assurance will likely gain rapid interest.

## Publication and Adoption Opportunities

**Viable Publication Venues:** Both frameworks sit at the intersection of AI, security, and enterprise IT, which opens multiple avenues for publication:

- *Academic and Research:* ArXiv is an immediate option for disseminating these framework ideas to a broad audience; indeed, similar work on zero trust agent identity has appeared there <sup>19</sup>. For EAAGF, an academic-style whitepaper could be released on arXiv or presented at conferences like **AAAI/ACM AI Ethics and Society**, **IEEE International Conference on AI** (for governance), or workshops at **NeurIPS** or **ICLR** that focus on safe and trustworthy AI. The Zero Trust framework, with its security slant, could be submitted to **IEEE Security & Privacy (Oakland)** workshops, **ACSAC** or **NDSS** workshops on AI security, or presented at **NIST** convenings (NIST has shown interest by issuing an RFI on AI agent security <sup>34</sup>).
- *Industry Conferences:* For EAAGF, industry forums like **Gartner Security & Risk Summit**, **OECD AI Governance events**, or **World Economic Forum** (which often issues whitepapers on AI governance)

would be ideal to launch a whitepaper. The framework could also be introduced at trade conferences such as **AI Summit (Governance track)** or sector-specific events (e.g., a banking risk management conference or **RSAC (RSA Conference) 2026** where AI governance is in focus). The Zero Trust agent framework would fit well as a talk or whitepaper at **RSA Conference**, **Black Hat**, or **DEF CON** (to reach security practitioners), as well as cloud security conferences (the **Cloud Security Alliance Summit**, **Identiverse** for digital identity, etc.). Even cross-domain conferences like **SXSW** or **Web Summit** might be interested, given the buzz around AI agents.

- *Whitepapers and Standards:* Both frameworks could be published as **whitepapers or technical reports** through recognized bodies. For instance, EAAGF could be released as an **enterprise AI governance guide** by a consortium of companies or an open industry group (with press coverage in trade media). The Zero Trust framework might be proposed as an extension to existing **Zero Trust Architecture** documentation or through the **CSA** (which could consider it under their Zero Trust or AI Safety initiative). Additionally, standardization bodies like the IEEE or ISO might be interested if positioned as a starting point for standards (e.g., an IEEE workshop on AI agent governance).
- *GitHub/Open-Source Repositories:* To encourage practical adoption, especially for the Zero Trust framework, publishing reference implementations or tools on GitHub is valuable. For example, a repository containing templates for agent credentials, sample policy engines, or scripts for monitoring agent actions would attract developers and DevOps teams. This also invites open-source communities (perhaps those around **LangChain**, **AutoGPT**, etc.) to contribute and use the frameworks.
- *Trade Press and Media:* To gain mindshare, the concepts should be introduced in accessible terms via tech media. Outlets like **MIT Technology Review**, **Wired**, **VentureBeat**, **TechCrunch** or enterprise IT publications (**CIO.com**, **Dark Reading**, **InfoWorld**) could feature op-eds or exclusive announcements. For instance, a piece titled "**Why Enterprises Need an AI Agent Governance Framework**" in Harvard Business Review or Forbes would target executives and underscore EAAGF's relevance. Likewise, the zero trust agent framework could be pitched to cybersecurity magazines (CSO Online, SC Magazine) as the answer to rising AI security fears. Early coverage will establish the frameworks' names in the industry lexicon.

### **Communities and Early Adopters:**

- EAAGF is likely to gain traction fastest in **highly regulated industries** (finance, healthcare, government) where strict governance is non-negotiable. These sectors have active communities (e.g., bank consortiums on AI risk, healthcare AI ethics groups) that would welcome a clear governance model. Government agencies or international organizations focusing on AI policy (like the EU's AI regulatory community or the US's NIST and NTIA in AI governance) might also pilot such a framework. Engaging with groups like **FINRA** or **FDA's AI initiatives**, the **OECD AI Policy Observatory**, or national AI governance pilot programs could fast-track adoption. Academic circles in AI ethics/governance and think tanks (e.g., **Partnership on AI**) are another community where EAAGF can gain intellectual backing.
- The Zero Trust Authentication Framework will resonate with the **cybersecurity and identity management community**. Early adopters could include tech companies building multi-agent platforms (who need security solutions now), cloud providers, and any enterprise already embracing

zero trust for their human users (they will naturally extend it to AI agents). Communities like **OWASP** (possibly creating an “Agent Security Project”), the **Identity Defined Security Alliance (IDSA)**, and security-focused Slack/Discord groups could be tapped. Additionally, cross-organization collaborations – for example, a few companies creating a shared “trust network” for agents using this framework – could showcase its value in real-world B2B agent interactions (imagine multiple banks using a common standard to trust each other’s AI agents safely).

**Summary of Novelty and Overlap:** The following table summarizes each framework’s novelty, similar existing frameworks, and overlap with current efforts:

Framework	Similar Existing Frameworks & Efforts	Novelty & Differentiation	Overlap & Positioning
<b>Enterprise AI Agent Governance Framework (EAAGF)</b> <i>&lt;br&gt;(Governance, oversight, maturity model)</i>	<ul style="list-style-type: none"> <li>- <i>AI Governance Standards</i>: NIST AI RMF, ISO 42001 (general AI risk management, not agent-specific) <sup>4</sup></li> <li><sup>5</sup> .&lt;br&gt;- <i>Agentic AI Governance</i>: Industry concepts from Palo Alto Networks, etc., stressing delegated authority and runtime controls for agents <sup>1</sup> <sup>2</sup> .&lt;br&gt;- <i>Autonomy Levels</i>: CSA's proposed 6-level autonomy taxonomy for agent behavior; various "levels of AI autonomy" blogs paralleling the idea <sup>6</sup> .&lt;br&gt;- <i>Integrity Models</i>: Acuvity's Agent Integrity Framework (5-level maturity focusing on security controls) <sup>7</sup> .&lt;br&gt;- <i>Multi-Agent Oversight</i>: McKinsey and others urging multi-agent oversight, human-in-the-loop triggers, and governance processes for agent deployments <sup>8</sup> <sup>9</sup> .</li> </ul>	<p><b>- Integrative Maturity</b></p> <p><b>Model:</b> Combines policy, technical controls, and risk scoring into one progression model (filling a gap where current frameworks address pieces in isolation).&lt;br&gt;- <b>Layered Control Approach:</b> Explicit control layers (audit logs, supervision checkpoints, autonomy limits) for <b>operational governance</b> – beyond static policy documents, it provides a live governance "control plane" for AI agents (novel in enterprise IT contexts).&lt;br&gt;- <b>Risk Scoring Mechanism:</b> Introduces quantifiable risk metrics for agent actions to dynamically adjust oversight – an innovation compared to qualitative risk categories in regulations. <sup>3</sup> &lt;br&gt;- <b>Enterprise Focus:</b> Tailored for regulated industries' needs (compliance, auditability, accountability) in deploying AI agents, whereas most academic frameworks are generic.</p>	<p><b>- Overlap:</b> Shares goals with AI risk frameworks (ensuring safe, compliant AI use) and with agent security efforts. Partial overlap with CSA and Acuvity on autonomy and security controls, indicating validation of the need.&lt;br&gt;- <b>Positioning:</b> Differentiate by breadth: EAAGF covers <b>full lifecycle governance</b> (from design to deployment to monitoring) with a maturity roadmap, whereas others are either high-level (NIST) or narrow (security-only). It can be positioned as a superset that <b>aligns with</b> standards (e.g., map EAAGF levels to NIST/ISO requirements <sup>35</sup>) but adds agent-specific depth. Emphasize real-world practicality – it's a "<i>ready-to-implement</i>" framework for CIOs/CISOs, not just theory. Collaborating with standards bodies or showcasing case studies (pilot implementations in a bank, etc.) will strengthen its unique value.</p>

---

<p><b>Zero Trust Authentication Framework for Autonomous Agents</b></p> <p>&lt;br&gt;(Identity, access, trust for AI agents)</p>	<ul style="list-style-type: none"> <li>- <i>Zero Trust Architecture:</i> Established cybersecurity approach (e.g., NIST SP 800-207) now being extended to non-human entities.</li> <li>- <i>CSA Agent Identity Framework:</i> CSA's proposals for agent IAM using DIDs/VCs and continuous verification <sup>17</sup> <sup>16</sup>; related arXiv paper by industry researchers defining a decentralized agent identity system <sup>19</sup>.</li> <li>- <i>Identity Management Solutions:</i> Okta's and others' guidelines on non-human identities (service accounts, bots) applying zero trust (short-lived credentials, IAM integration) <sup>20</sup>.</li> <li>- <i>Agent Communication Protocols:</i> Emerging standards (Anthropic, Google, IBM) for agent interactions, which implicitly need zero-trust principles (authentication, least privilege) <sup>27</sup>.</li> <li>- <i>AI Security Research:</i> Security community identifying agent attack vectors (spoofed identities, cross-agent exploits) and calling for frameworks (e.g., OWASP, Keyfactor)</li> </ul>	<p><b>- Novel IAM Paradigm:</b> Redefines digital identity for AI agents (using cryptographic proofs of identity and intent) – moving beyond simplistic API keys to rich <b>Agent IDs</b> with provenance and context (very new concept) <sup>19</sup>.</p> <p><b>- Dynamic Trust &amp; Authorization:</b> Implements real-time, fine-grained authorization that adapts to agent behavior and environment ("trust no one, verify always" applied continuously). While zero trust for humans exists, doing so for autonomous decision loops (with agents potentially chaining actions) is innovative.</p> <p><b>Federated Cross-Domain Trust:</b> Proposes a method for agents from different organizations to trust each other's credentials and behavior safely – essentially creating an <b>inter-organizational trust fabric</b> for AI agents, which is unique. Leverages cutting-edge tech like decentralized credentials and possibly blockchain or PKI in new ways for AI.</p> <p><b>Holistic Security Coverage:</b> Merges identity security with AI behavior monitoring (intent alignment)</p>	<p><b>- Overlap:</b> Strong alignment with CSA's zero-trust IAM vision and academic proposals (indicating the idea is timely). Overlaps with identity vendors' approaches to some extent (they handle machine identities but may not have full behavioral trust scoring). The concept of continuous agent verification is being discussed in pieces, validating the need.</p> <p><b>- Positioning:</b> Emphasize completeness: whereas others may offer parts (just identity management or just monitoring), this framework is an <b>end-to-end reference model</b>. It can be the blueprint that ties together identity, access, and behavior into one trust framework for agents. Position it as complementary to existing zero trust efforts: organizations likely have zero trust for users – now here's the counterpart for AI agents. By presenting it as a <i>standard or open framework</i>, we invite collaboration rather than competition. Highlight any areas not covered by others – for instance, the <b>Agent Name Service and global policy layer</b> (from the CSA model) can be a selling point as a novel infrastructure component to manage agent identities globally <sup>16</sup>. Additionally,</p>
--	---	---	--

Framework	Similar Existing Frameworks & Efforts	Novelty & Differentiation	Overlap & Positioning
	blogs) to secure agent interactions <sup>36</sup> <sup>37</sup> .	checks, anomaly detection) to cover threats classic security tools miss (like the logic-layer attacks described in CSA's work) <sup>17</sup> <sup>32</sup> .	stress <b>future readiness</b> : as agent ecosystems grow, a unified zero trust framework will save companies from patchwork solutions. Early adoption positions organizations ahead of emerging regulations (e.g., EU AI Act's requirements for oversight and audit trails for autonomous systems <sup>38</sup> <sup>39</sup> ).

## Conclusion

Both the Enterprise AI Agent Governance Framework and the Zero Trust Authentication Framework for Autonomous Agents address urgent gaps in how organizations manage the next generation of AI systems. Existing literature and industry activity show that these ideas are timely – there are partial precedents and parallel efforts – but each framework offers a **novel integration of concepts** that could set it apart.

EAAGF can become a cornerstone for **responsible AI deployment in enterprises**, translating high-level principles into an actionable program that evolves as organizations mature. In a landscape where only a small fraction of firms feel prepared for autonomous AI <sup>40</sup>, EAAGF could quickly gain support as companies and regulators alike seek concrete governance solutions. By aligning with known frameworks (NIST, ISO, EU AI Act) while introducing agent-specific controls, it can be marketed as the practical way to achieve compliance and confidence in AI agent deployments <sup>41</sup> <sup>42</sup>.

The Zero Trust agent framework, on the other hand, positions itself at the forefront of **AI security innovation**. As enterprises realize that AI agents are essentially new “users” on their networks (often with elevated powers), the zero trust approach will become non-negotiable. This framework’s focus on verifiable identity and continuous trust will help organizations answer the critical question posed by security experts: *“When an AI agent acts on your behalf, how do you know it is working for you – and not subverted or misused?”* <sup>43</sup> <sup>44</sup>. By pre-empting standards and offering a blueprint, it can influence how the industry builds secure agent ecosystems from the ground up.

**Next Steps:** To capitalize on their novelty, the proponents of these frameworks should consider early publication (e.g., arXiv or a well-publicized whitepaper) to stake thought leadership. Engaging with communities (CSA working groups, OWASP, IEEE, etc.) will refine the ideas with peer input and build credibility. Pilot projects or simulations demonstrating the frameworks in action (for example, a sandboxed enterprise agent with EAAGF controls enabled vs. disabled, or an inter-company agent transaction using the zero trust framework) would provide compelling evidence of value.

In positioning these frameworks, it will be important to acknowledge existing work and position the frameworks as **collaborative evolutions** rather than reinventing the wheel. For instance, citing how EAAGF complements the NIST AI RMF by adding an agent maturity dimension, or how the zero trust framework builds on CSA's guidelines by providing a deployable reference architecture, will show that these ideas stand on the shoulders of prior art while charting new territory.

Ultimately, the rapid pace of AI agent adoption means enterprises and the security/governance industry are hungry for solutions. By being comprehensive yet pragmatic, these frameworks can quickly gain traction as go-to references – much like earlier frameworks did in their domains (e.g., ITIL for IT service management or Zero Trust for network security). The comparative analysis suggests strong interest and some competition in each area, but also that there is ample room for a well-crafted framework to lead the conversation. With strategic outreach and perhaps open-source elements, EAAGF and the Zero Trust Authentication Framework could become foundational in enabling **safe, trustworthy, and governed AI agent deployments** in the coming years.

---

1 2 3 A Complete Guide to Agentic AI Governance - Palo Alto Networks

<https://www.paloaltonetworks.com/cyberpedia/what-is-agentic-ai-governance>

4 5 8 35 41 42 AI Governance Frameworks & Best Practices for Enterprises 2026

<https://onereach.ai/blog/ai-governance-frameworks-best-practices/>

6 Autonomy Levels for Agentic AI | CSA

<https://cloudsecurityalliance.org/blog/2026/01/28/levels-of-autonomy>

7 23 24 25 26 31 32 34 43 44 The Agent Integrity Framework: The New Standard for Securing Autonomous AI - Acuity

<https://acuity.ai/the-agent-integrity-framework-the-new-standard-for-securing-autonomous-ai/>

9 10 27 36 37 40 Agentic AI security: Risks & governance for enterprises | McKinsey

<https://www.mckinsey.com/capabilities/risk-and-resilience/our-insights/deploying-agentic-ai-with-safety-and-security-a-playbook-for-technology-leaders>

11 12 13 14 15 20 21 38 39 Agentic AI Frameworks: Identity, Security, Governance | Okta

<https://www.okta.com/identity-101/agentic-ai-framework/>

16 17 18 29 30 Fortifying the Agentic Web: A Unified Zero-Trust Architecture for AI

<https://cloudsecurityalliance.org/blog/2025/09/12/fortifying-the-agentic-web-a-unified-zero-trust-architecture-against-logic-layer-threats>

19 28 A Novel Zero-Trust Identity Framework for Agentic AI: Decentralized Authentication and Fine-Grained Access Control

<https://arxiv.org/html/2505.19301v2>

22 Zero trust for agentic systems: Managing non-human identities at scale

<https://www.hashicorp.com/blog/zero-trust-for-agentic-systems-managing-non-human-identities-at-scale>

33 Zero Trust for Agentic AI Security

<https://www.keyfactor.com/solutions/secure-ai-agents/>