Agent Overlay Whitepaper

Trust, Security & Compliance for the Agentic Workforce

Author: Andor Kesselman
Organization: Agent Overlay, Inc
Version: v0.0.2

Date: May 12, 2025

Description:

A framework for trust, security, and compliance in the agentic workforce.

Contents

1	Introduction	2
2	Foundational Concepts	2
3	Problem Statement	3
	3.1 Increasing Complexity of the Agentic Workforce	3
	3.2 Agents and the Legal Landscape: Tools for Better Governance	
	3.3 Compliance Risk and the Cost of Governance	
	3.4 Control Is Moving Away: The Need for Decentralized Governance Tools	
	3.5 Unit Economics and Liability: The Critical Barrier to Scale	
	3.5.1 Parameters	
	3.5.2 Expected Net Benefit	
	3.5.3 Economic Viability Condition	
	3.5.4 Need for Better Governance Tools	
	3.6 What This Means For the Agentic Workforce	
4	Agent Overlay Overview	6
	4.1 How Decentralized Protocols Extend the Control Plane Across Business Contexts	6
	4.2 Agent Overlay High Level Architecture	6
	4.3 Key Product Positions	
	4.3.1 Neutrality	
	4.3.2 Product Focus	
	4.3.3 Scalability	
5	Conclusion	8

1 Introduction

The rapid advancement of large language models (LLMs) has sparked an unprecedented wave of innovation in AI-driven processes. While artificial intelligence (AI) and machine learning (ML) are not new fields, the speed of adoption, particularly in business contexts, has exploded, driven by the promise of immediate problem-solving and operational efficiency. Organizations now leverage AI to automate workflows, enhance decision-making, and unlock new business opportunities.

However, the complexity of modern business environments underscores the limitations of generalized AI solutions. No single agent can address all problems effectively. Instead, the future lies in custom, domain-specific agents tailored to specific industries or workflows. These agents, designed with specialized knowledge and capabilities, offer precision and adaptability that general-purpose AI lacks.

This shift toward specialized agents brings new challenges. As businesses integrate AI into critical processes, they encounter risks related to data privacy, governance, auditability, and operational control. Ensuring trust in these systems demands an overlay of governance mechanisms that can monitor, validate, and guide AI operations. Such mechanisms must balance the need for control with the flexibility to innovate, emphasizing decentralization, minimal oversight, and ease of use.

For AI to fulfill its potential in business, it must not only solve problems but also inspire trust. In many cases, governance frameworks must ensure transparency and accountability while remaining accessible to businesses of all sizes.

This paper proposes a novel approach to address these challenges. It outlines a decentralized control framework that leverages decentralized identifiers, agent delegation, event chains, and context overlays to create a robust network of agent-driven processes. This solution provides a flexible and scalable framework for constructing complex agent workflows without relying on a centralized workflow management system that is not viable for complex organizational operation chains.

2 Foundational Concepts

- Agent Identifiers Agent identifiers are critical components in any multi-agent system, serving as unique references for individual agents within the system. These identifiers enable agents to interact, collaborate, and authenticate their actions. Traditional systems rely on centralized identifiers, which often introduce bottlenecks and vulnerabilities due to dependency on a single point of control. Moving toward decentralized approaches improves scalability, resilience, and security in agent interactions.
- Decentralized Identifiers (DIDs) Decentralized Identifiers (DIDs) extend the concept of agent identifiers by providing a self-sovereign and cryptographically secure framework for identification. Unlike centralized identifiers, DIDs are managed directly by their subjects, eliminating reliance on centralized authorities. In multi-agent systems, DIDs facilitate secure communication and verifiable interactions. Agents can use DIDs to establish trust, authenticate themselves, and record interactions transparently. These features are foundational to enabling decentralized governance and ensuring system-wide accountability.
- Delegation Chains Delegation chains are dynamic pathways through which tasks, responsibilities, and permissions flow between agents. In decentralized systems, these chains allow agents to delegate tasks to other agents based on real-time evaluations of capabilities and contexts. Effective delegation chains ensure that the most qualified agent executes each task, optimizing system performance. They also incorporate mechanisms for grant negotiation, enabling agents to formalize and revoke permissions dynamically, thus maintaining flexibility and control.
- Verifiable Event Chains Verifiable event chains document the sequence of actions and decisions made by agents within a system. These chains create an immutable audit trail that can be used for compliance, accountability, and troubleshooting. By leveraging cryptographic methods, verifiable event chains ensure that each recorded event is tamper-proof and attributable to a specific agent. In decentralized overlays, these chains are critical for maintaining transparency and trust, as they provide stakeholders with a clear view of system operations.

- Context Overlays Context overlays provide additional layers of metadata and situational awareness that agents use to make informed decisions. These overlays enable agents to understand the broader environment in which they operate, including system policies, resource availability, and the goals of other agents. Context overlays also help in resolving conflicts and prioritizing tasks by providing agents with a shared understanding of their operational context. In decentralized systems, context overlays are instrumental in aligning agent actions with system-wide objectives while maintaining autonomy.
- Agent Discovery Agent discovery is the process through which agents in a multiagent system locate and identify other agents for collaboration, delegation, or communication.
- Model Context Protocol (MCP) The Model Context Protocol (MCP) defines the framework through which agents understand, share, and operate within a shared context. It serves as a foundation for aligning agent actions, ensuring coherent decision-making, and facilitating interoperability in multi-agent systems. MCP enables agents to exchange relevant contextual data while maintaining autonomy and adaptability.
- Machine Contract Negotiation Protocol (MCNP) The Machine Contract Negotiation Protocol (MCNP) ensures that before an agent executes any task or interaction, a contract must be proposed, reviewed, and accepted. This protocol governs the lifecycle of contractual agreements between AI agents and human operators, ensuring that terms are explicitly agreed upon before execution.

3 Problem Statement

3.1 Increasing Complexity of the Agentic Workforce

The proliferation of specialized AI agents will led to a highly fragmented digital ecosystem. Instead of relying on a single centralized system, modern organizations deploy a diverse array of narrowly focused agents tailored to specific functions such as scheduling, finance, HR, and customer support. This shift necessitates varied hosting and management strategies depending on the use case. In industries where Private AI is a critical requirement, organizations face a choice similar to that seen in cryptocurrency custody models: some opt for privately hosted (non-custodial) AI solutions to retain full control over their data and operations, while others leverage third-party (custodial) AI services for ease of integration and reduced infrastructure complexity.

This fragmentation means that effective control cannot assume a centralized model. Instead, it demands more sophisticated, decentralized control mechanics. The challenge is to seamlessly integrate these specialized agents so that, despite their autonomy, the overall system remains efficient, secure, and aligned with organizational objectives.

3.2 Agents and the Legal Landscape: Tools for Better Governance

UETA is a legal framework that defines how transactions involving agents are handled, setting standards for electronic contracts, signatures, and clear attribution of actions. In contrast, GDPR and other data governance regulations focus on protecting personal and sensitive information by imposing strict data privacy and security requirements.

The increasing complexity of these distinct legal and regulatory environments highlights the urgent need for sophisticated tools to manage contractual engagements with agentic systems. On one hand, UETA ensures that AI agents can participate in legally binding transactions by enforcing clear accountability and structured interactions. On the other hand, GDPR demands rigorous data protection measures, ensuring that any data processed by these systems complies with high privacy standards.

To navigate these parallel challenges, organizations must develop advanced compliance solutions that not only streamline contract enforcement and accountability for AI agents but also integrate robust data governance controls. Such tools would enable real-time monitoring, auditing, and control

over AI activities, ensuring that all interactions adhere to the distinct requirements of both transactional law and data privacy regulations. This dual approach is essential for safely harnessing the transformative potential of AI in a regulated digital landscape.e fully realized in practice.

3.3 Compliance Risk and the Cost of Governance

As AI agents multiply, so does the complexity and cost of compliance. For basic AI compliance in the EU compliance expenses can represent up to 17% of total AI investments Data Innovation Report 2021. By 2025, the Artificial Intelligence Act is expected to cost the EU \$30B alone. In some deployment projects, companies have reported compliance costs exceeding \$340,000—often dwarfing their R&D spending Harvard Student Review. Moreover, penalties for non-compliance have been severe, with fines in the hundreds of millions reported in cases like Didi and Amazon Holistic AI. These figures highlight that without scalable compliance and governance mechanisms, organizations will be reluctant to embrace AI systems, as even a single security or regulatory breach can result in substantial financial losses.

3.4 Control Is Moving Away: The Need for Decentralized Governance Tools

The landscape of AI is shifting from centralized control to a distributed network of specialized agents. In Europe, examples such as the eiDAS Trust Framework and elements of the evolving EU Artificial Intelligence Act underscore the necessity for decentralized yet coherent control mechanisms. With disparate agents handling different functions—from scheduling to finance—the absence of an overarching, scalable control system leads to risks of misalignment, fragmented accountability, and increased liability.

For instance, if a budgeting agent and a scheduling agent are developed by different vendors with incompatible control mechanisms, an organization may face significant operational disruptions. This illustrates why new systems must be engineered to manage these decentralized relationships—ensuring that every agent adheres to rigorous standards of transparency, error correction, and accountability. These standards are not provided by UETA or eiDAS alone, but they offer valuable precedent that must be built upon with modern tools.

3.5 Unit Economics and Liability: The Critical Barrier to Scale

The economic viability of an agentic workforce hinges on a fundamental trade-off between the operational benefits and the risks of liability. In high-stakes environments, the cost of a mistake can be catastrophic. For example, it is not feasible for Tesla to roll out Full Self-Driving (FSD) technology widely unless the liability exposure is drastically reduced, either technically or legally. Similarly, if AI agents in business transactions expose organizations to multi-million-dollar fines or legal uncertainties, the promise of enhanced efficiency and productivity will remain unrealized.

Until liability and compliance costs are significantly lower than the operational upside, organizations will hesitate to deploy agentic systems at scale—especially in areas involving sensitive data and critical operations.

The economic viability of an agentic workforce hinges on the trade-off between operational benefits and the risks of liability and compliance. Below, we define key parameters and construct a mathematical model to analyze the conditions under which AI-driven agents can be deployed profitably.

3.5.1 Parameters

- **B** Per-Unit Operational Benefit The additional value (in dollars) generated by one unit of AI-agent output compared to a non-AI alternative.
- C_O Per-Unit Operational Cost The cost of running and maintaining the AI agent for each unit of output (e.g., computing, development, and infrastructure costs).
- N Number of Units The total number of tasks, transactions, or actions performed by the AI agent over a given period.
- p Probability of Error or Negative Event The likelihood that an AI agent's action results in a costly mistake, compliance violation, or legal liability.
- L Liability Cost per Error The expected monetary damage, fine, or other liability associated with each negative event.
- C_{compliance} Compliance & Risk Management Cost The total fixed cost of meeting regulatory, auditing, and legal requirements.
- I Initial Implementation Cost A one-time investment required to deploy the AI agent system, including system integration and legal setup.

3.5.2 Expected Net Benefit

The expected net benefit (NB) of deploying AI agents is:

$$NB = (B \times N) - (I + C_O \times N) - (p \times L \times N) - C_{\text{compliance}}$$

Where:

- 1. $B \times N$ represents the total operational gains.
- 2. $I + C_O \times N$ represents implementation and operational costs.
- 3. $p \times L \times N$ captures the expected liability from AI errors.
- 4. $C_{\text{compliance}}$ is the fixed cost of ensuring regulatory compliance.

3.5.3 Economic Viability Condition

For AI deployment to be viable, the net benefit must be positive:

$$B \times N - (I + C_O \times N) - (p \times L \times N) - C_{\text{compliance}} > 0$$

Rearranging:

$$N \times (B - C_O - p \times L) > I + C_{\text{compliance}}$$

This equation shows that for large-scale deployment, the **per-unit net margin** $B - C_O - p \times L$ must be positive and large enough to offset the fixed compliance and setup costs.

1. Interpretation

High-Stakes Environments (e.g., Autonomous Driving)

If p (probability of failure) and L (liability cost) are high, $p \times L$ can outweigh operational benefits, discouraging deployment. Tesla, for example, cannot deploy Full Self-Driving (FSD) unless it drastically reduces liability risks.

Regulatory-Heavy Sectors (e.g., Healthcare, Finance)

Compliance costs ($C_{\text{compliance}}$) are high, requiring advanced legal safeguards, increasing the threshold at which AI adoption becomes viable.

Data-Intensive Scenarios (e.g., GDPR, HIPAA Compliance)

Data governance laws create additional compliance overhead that may not fit into traditional liability models but contribute significantly to overall costs.

3.5.4 Need for Better Governance Tools

Given the complexities of liability, compliance, and operational costs, organizations need automated contract enforcement, real-time risk management, and AI governance frameworks to lower p, minimize L, and streamline $C_{\text{compliance}}$. Without such tools, organizations will hesitate to deploy AI agents in legally sensitive domains.

3.6 What This Means For the Agentic Workforce

The increasing fragmentation of the agentic workforce, driven by specialization, demands innovative and decentralized control and compliance solutions. The legal frameworks such as UETA and the eiDAS Trust Framework provide essential starting points and legal precedents for trusted electronic interactions. However, these frameworks are not complete solutions for the dynamic challenges of today's decentralized, specialized agents. Modern tools are needed to dynamically enforce legal principles, ensure transparent auditing, and manage risk across diverse AI systems.

This is precisely why a solution like **Agent Overlay** is essential: it offers a decentralized control overlay that unifies governance, reduces liability, and unlocks the full potential of an efficient and trustworthy AI workforce. Without such a framework, organizations will continue to face prohibitive compliance costs and legal uncertainties that hinder the transformative potential of agentic AI.

4 Agent Overlay Overview

Agent Overlay is a comprehensive platform of software, protocols, and infrastructure designed to secure and govern the increasingly fragmented agentic workforce. As a Multi-Agent Decentralized Control Overlay (Agent Overlay) for Business Security & Governance, it empowers organizations to conduct secure, cross-ecosystem interactions that transcend traditional business boundaries.

- For Developers, Agent Overlay offers a robust SDK that streamlines integration with the network overlay and unlocks powerful features in both new and existing applications.
- For **Operators**, Agent Overlay provides a clear, intuitive portal for real-time oversight, policy configuration, and continuous management of the agentic control layer.

These tools work together to deliver a secure, compliant, and efficient environment. Specialized AI agents can operate autonomously while still upholding organizational objectives and meeting all relevant legal and regulatory standards.

4.1 How Decentralized Protocols Extend the Control Plane Across Business Contexts

Agent Overlay's decentralized protocols form the foundation of its cross-boundary control plane. By distributing authority and verification among multiple stakeholders instead of relying on a single, centralized system, organizations of any size or structure can adopt a shared governance model to regulate AI agents. The result is a flexible ecosystem with robust security and compliance capabilities, built on protocols that unify auditing, reporting, and policy enforcement across diverse operational contexts, including:

- Identity Distributed identifiers (DIDs) ensure each agent has a cryptographically secure, portable identity. This design mitigates single points of failure in authentication and strengthens overall system resilience.
- Delegation Protocols Provide a mechanism to grant and revoke specific tasks, permissions, or capabilities across autonomous agents in real time. This granular delegation system ensures that each agent only operates within the bounds set by the delegator, striking a balance between autonomy and controlled authority.

- Machine Contract Negotiation Protocol (MNCP) Employs targeted rule sets and real-time frameworks to govern agent behavior, helping organizations align with critical statutory requirements such as GDPR, HIPAA, or emerging AI-specific regulations. By requiring explicit acceptance of contractual terms before an agent proceeds, MNCP enforces compliance and traceability throughout interactions.
- Auditing and Event-Chaining Generates tamper-evident logs of agent actions and decisions, ensuring privacy and security without reliance on a public blockchain.

 These logs can be scrutinized for post-incident forensics or proactive compliance measures, offering a transparent but controlled view of agent activities.
- Cross-Boundary Collaboration Allows agents from different organizations or ecosystems to negotiate tasks, share data, and exchange capabilities securely. This approach preserves the ownership of proprietary assets while enabling seamless coordination, even across traditionally siloed environments.

4.2 Agent Overlay High Level Architecture

The Agent Overlay architecture comprises multiple layers to simplify security and compliance while allowing heterogeneous AI agents to cooperate within and across organizational silos:

- Network Overlay Layer: Manages decentralized identity, discovery, and secure communication channels for agent interactions.
- Governance & Policy Layer: Enforces consistent rules for delegation, auditing, and compliance. Operators can set policies that automatically apply to both internal and cross-ecosystem workflows.
- Agent Execution Layer: AI agents live here, interacting with the overlay to retrieve permissions, report their activities, and exchange capabilities with other agents.
- Developer & Operator Tools: Provide ready-to-use SDKs for building agentic applications and a user-friendly control panel for monitoring, policy configuration, and real-time management.

Through these layers, Agent Overlay supplies the trust, security, and compliance needed for

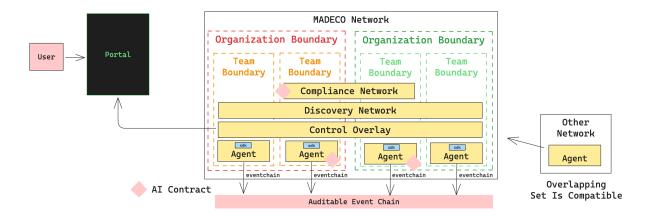


Figure 1: Agent Overlay High-Level Architecture

the modern agentic workforce—ensuring that AI adoption can scale confidently in even the most stringent regulatory or enterprise environments.

4.3 Key Product Positions

4.3.1 Neutrality

For AI-driven agentic systems to achieve widespread adoption, **neutrality is essential**. Unlike existing AI governance models that often prioritize proprietary ecosystems, Agent Overlay is designed as an **open**, **model-agnostic** control overlay. It does not favor specific industries, regulatory frameworks, or commercial entities. Instead, it ensures:

- Cross-Ecosystem Interoperability Agent Overlay enables secure, decentralized agentic interactions across business boundaries
- Flexible Compliance and Governance
 It integrates seamlessly with diverse compliance frameworks without imposing a singular approach.
- User-Controlled Decision-Making Operators can enforce their own rules, policies, and delegation mechanisms.
- Model Independence Enterprises retain the ability to use their preferred AI models while maintaining security and auditability.

By upholding neutrality, Agent Overlay prevents vendor lock-in and **empowers organizations to retain control** over their AI agents while ensuring compliance and trust.

4.3.2 Product Focus

Agent Overlay is **not a general-purpose AI** solution. It is purpose-built to secure and gov-

ern agentic transactions. Unlike broad AI platforms, Agent Overlay is laser-focused on:

- Decentralized Agentic Delegation Using verifiable credential chains to define trust and responsibility.
- Verifiable Event Chains Ensuring provable, auditable AI behavior through tamper-resistant logs.
- Agentic Term Exchange for Compliance Facilitating structured governance via negotiable agent contracts.
- Operator Injection Enabling real-time operator oversight and intervention within AI workflows.

By concentrating exclusively on **security**, **compliance**, **and governance**, Agent Overlay ensures that **AI agents are verifiable**, **accountable**, **and legally enforceable**, even in high-risk environments.

4.3.3 Scalability

To support enterprise-scale, mission-critical AI deployments, Agent Overlay must scale across operational, technical, and regulatory dimensions:

- Decentralized Architecture A distributed control overlay prevents reliance on a single entity, ensuring fault tolerance and regulatory independence.
- Technical Scalability Agent Overlay supports high-volume AI transactions with low latency and high reliability.
- Regulatory Adaptability The framework dynamically adjusts to evolving legal and compliance requirements without requiring a fundamental redesign.

• Enterprise Expansion – AI agent deployments can grow from pilots to full-scale integrations while maintaining governance, oversight, and trust.

By prioritizing **decentralization**, **compliance**, **and scalability**, Agent Overlay ensures **long-term sustainability** in regulated industries and high-stakes environments.

5 Conclusion

Interested in learning more? Check out www.madesm.com for more details.