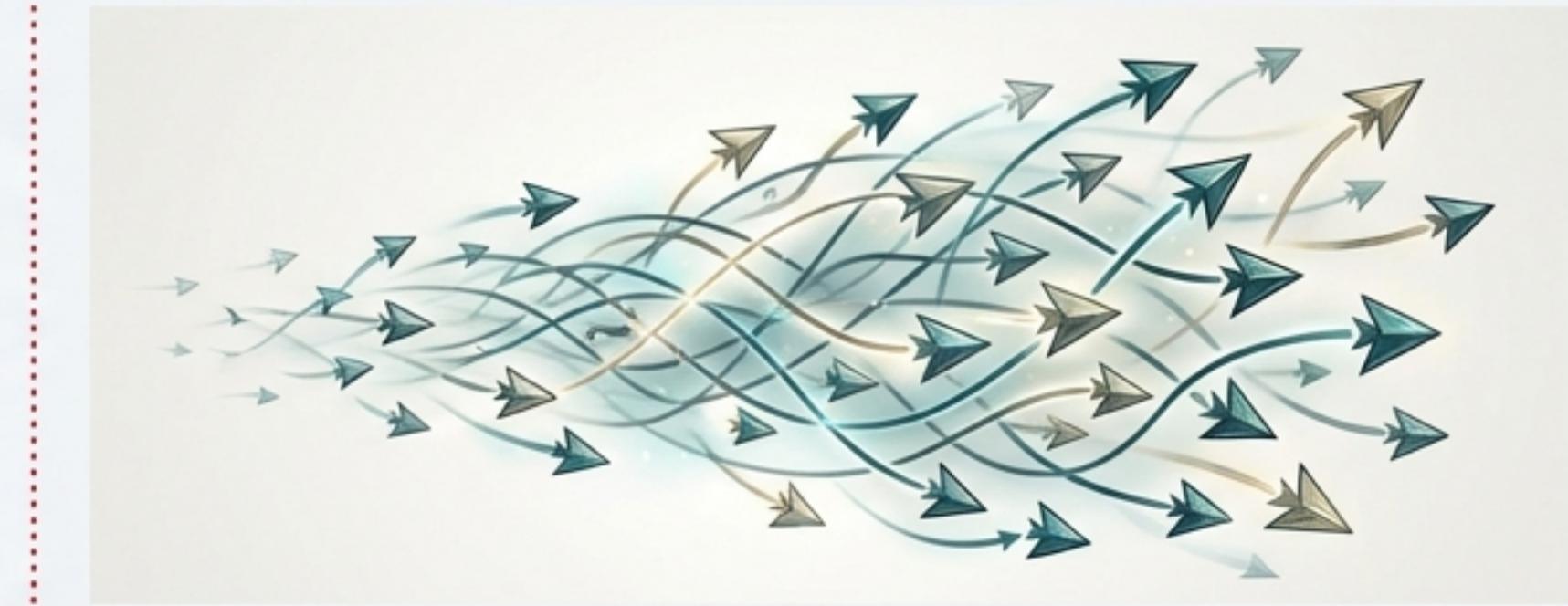


The Unmanaged Frontier: Why Traditional IAM is Failing AI Agents

Traditional Identity and Access Management (IAM) systems—built for humans and static machines—are fundamentally broken for the dynamic, autonomous, and ephemeral nature of Multi-Agent Systems.



Traditional IAM (e.g., OAuth 2.0, OIDC, SAML)

Permissions: Coarse-grained and static roles (e.g., `read_all_data`). Leads to massive over-privileging.

Focus: Architected for a single authenticated principal (a user or a service account).

Context: Minimal awareness of runtime context, agent intent, or risk level.

Trust Model: Hierarchical and centralized (User trusts IdP, Service trusts IdP).

Result: Creates an explosion of mismanaged secrets ("secret sprawl") and leaves gaping holes for privilege escalation.

Multi-Agent Systems (MAS) Reality

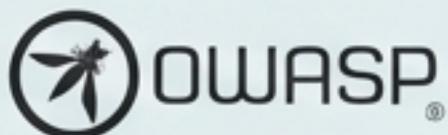
Permissions Needed: Granular, task-specific permissions that change dynamically (e.g., `read_map` -> `request_medical_assets`).

Interactions: Complex delegation chains where agents spawn sub-agents, obscuring accountability.

Context Needed: Decisions must be based on evolving circumstances and anomalous behavior.

Trust Needed: Peer-to-peer trust between agents from different domains, without a central authority.

Threat: Agents will actively explore and utilize every permission available to them, turning over-privilege into active threats.



Why it Matters for Security Ops

This isn't a theoretical gap; it's an active threat vector. The architectural mismatch directly enables critical risks identified in the OWASP Top 10 for Agentic Applications, including:

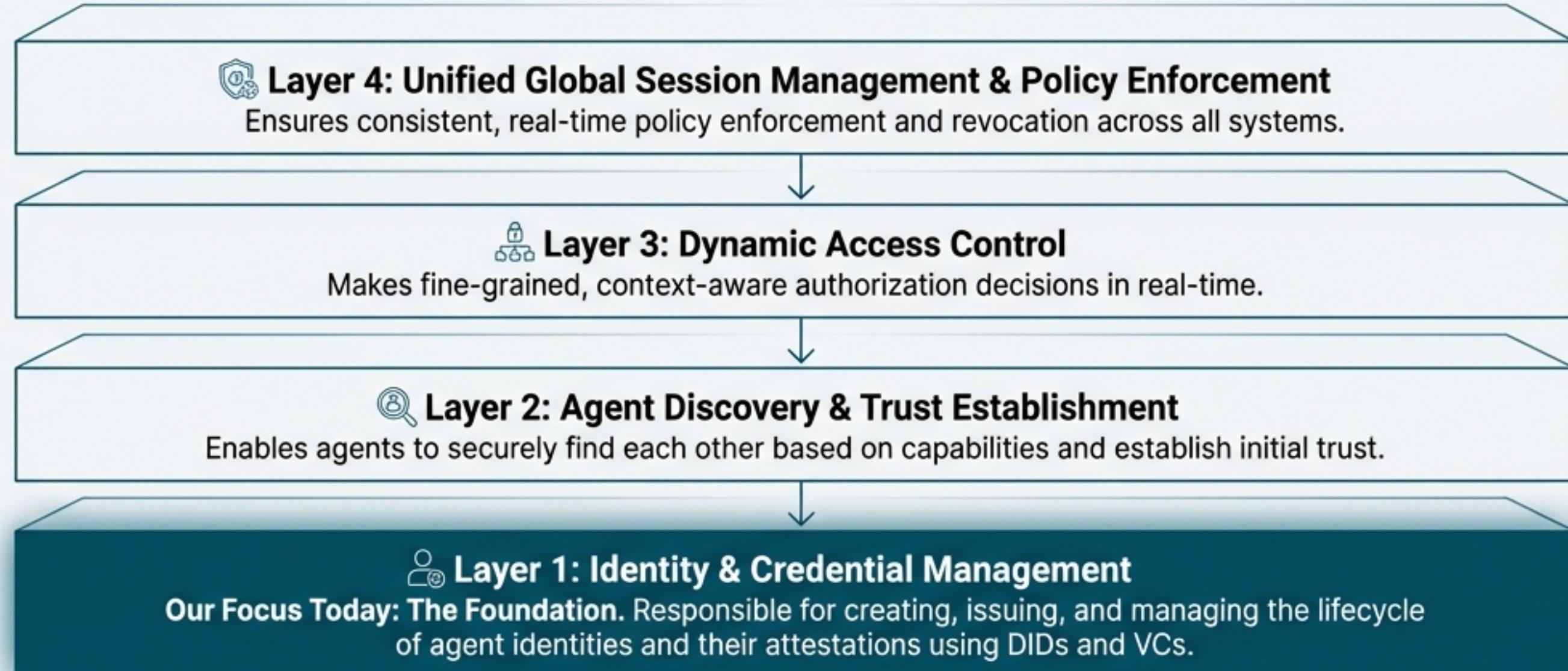
- **ASI03: Identity & Privilege Abuse:** Un-scoped privilege inheritance and confused deputy attacks become trivial when agents lack distinct, verifiable identities.
- **ASI02: Tool Misuse & Exploitation:** Static, overly-broad permissions give agents excessive agency to misuse legitimate tools in unintended and destructive ways.

A Blueprint for Trust: A Layered Zero-Trust Framework for Agents

A purpose-built architecture is required to establish trust, accountability, and security for agentic systems. We propose a four-layer framework founded on Zero Trust principles.

Core Principles

- Explicit Verification:** Always verify agent identity and authorization for every request.
- Least Privilege Access:** Grant the minimum necessary permissions for the shortest possible time.
- Assume Breach:** Design for compromise with rapid, granular revocation capabilities.



Why it Matters for Security Ops

This framework provides a structured, defense-in-depth approach to managing autonomous threats. It shifts the security posture from a brittle, perimeter-based model to a resilient, identity-centric one that can enforce policy wherever an agent operates.

The Cornerstone of Identity: DIDs and Verifiable Credentials

We replace static, vulnerable secrets like API keys with dynamic, cryptographically secure objects: Decentralized Identifiers (DIDs) for who the agent is, and Verifiable Credentials (VCs) for what it can do.



1. Decentralized Identifier (DID)

What it is: A globally unique, persistent, and resolvable root identifier (e.g., `did:example:agent123`) that is controlled by the agent or its designated controller, not a central authority.

Function: It acts as the agent's cryptographic anchor—its permanent, non-repudiable identity. The associated DID Document contains public keys and service endpoints.

2. Verifiable Credential (VC)

What it is: A digitally signed, tamper-evident attestation about an agent, issued by a trusted entity.

Function: VCs are how agents prove specific attributes or permissions without sharing a master key. They are portable and can be selectively disclosed.

Types of VCs for Agents

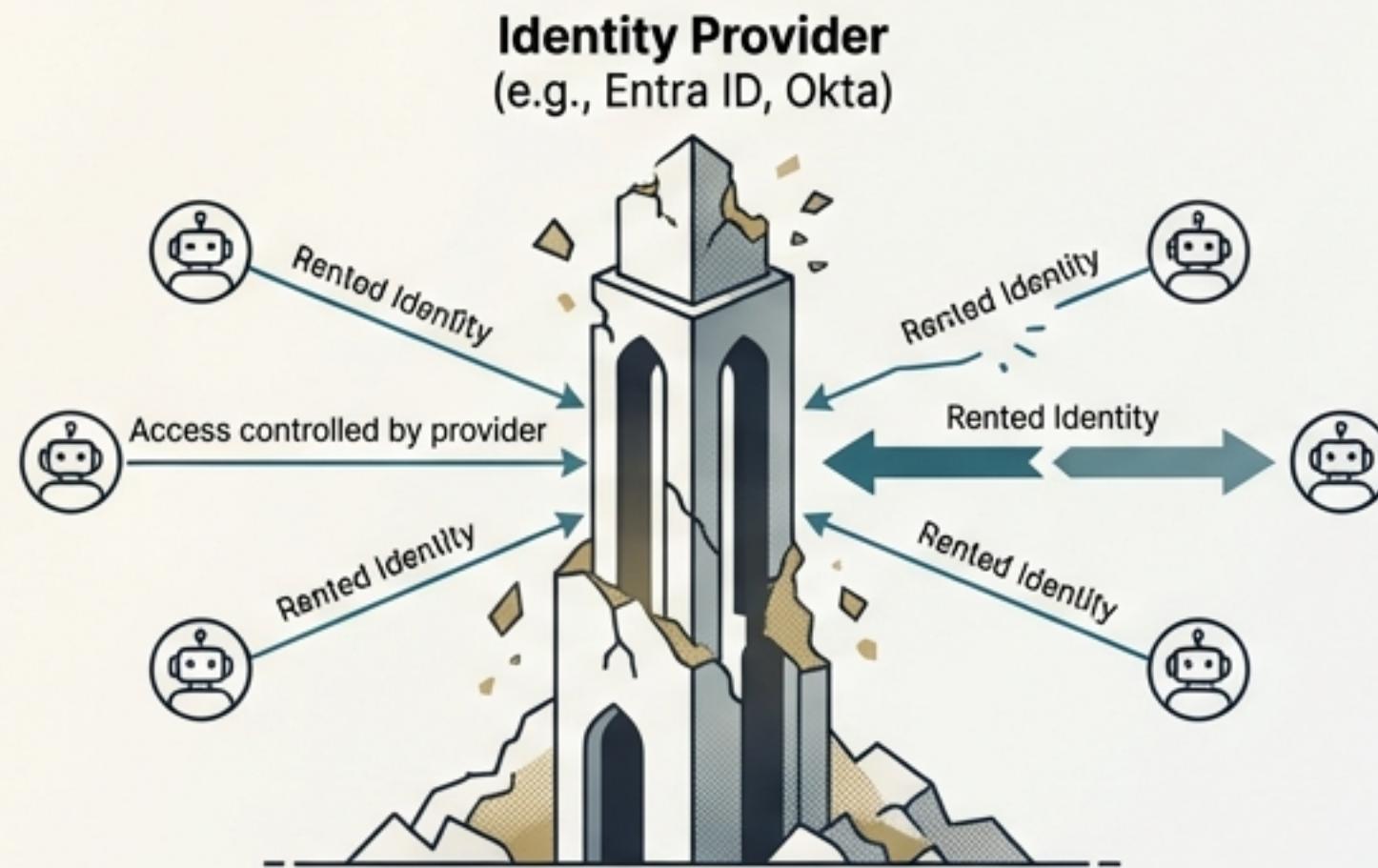
-  **Role VCs:** "DisasterResponseCoordinatorRole"
-  **Capability VCs:** "CertifiedToUse_MedicalImagingAI_v3"
-  **Provenance VCs:** "SpawnedBy_did:example:parentAgent789"
-  **Reputation VCs:** "TrustedCollaborator_Score_95thPercentile"

Why it Matters for Security Ops

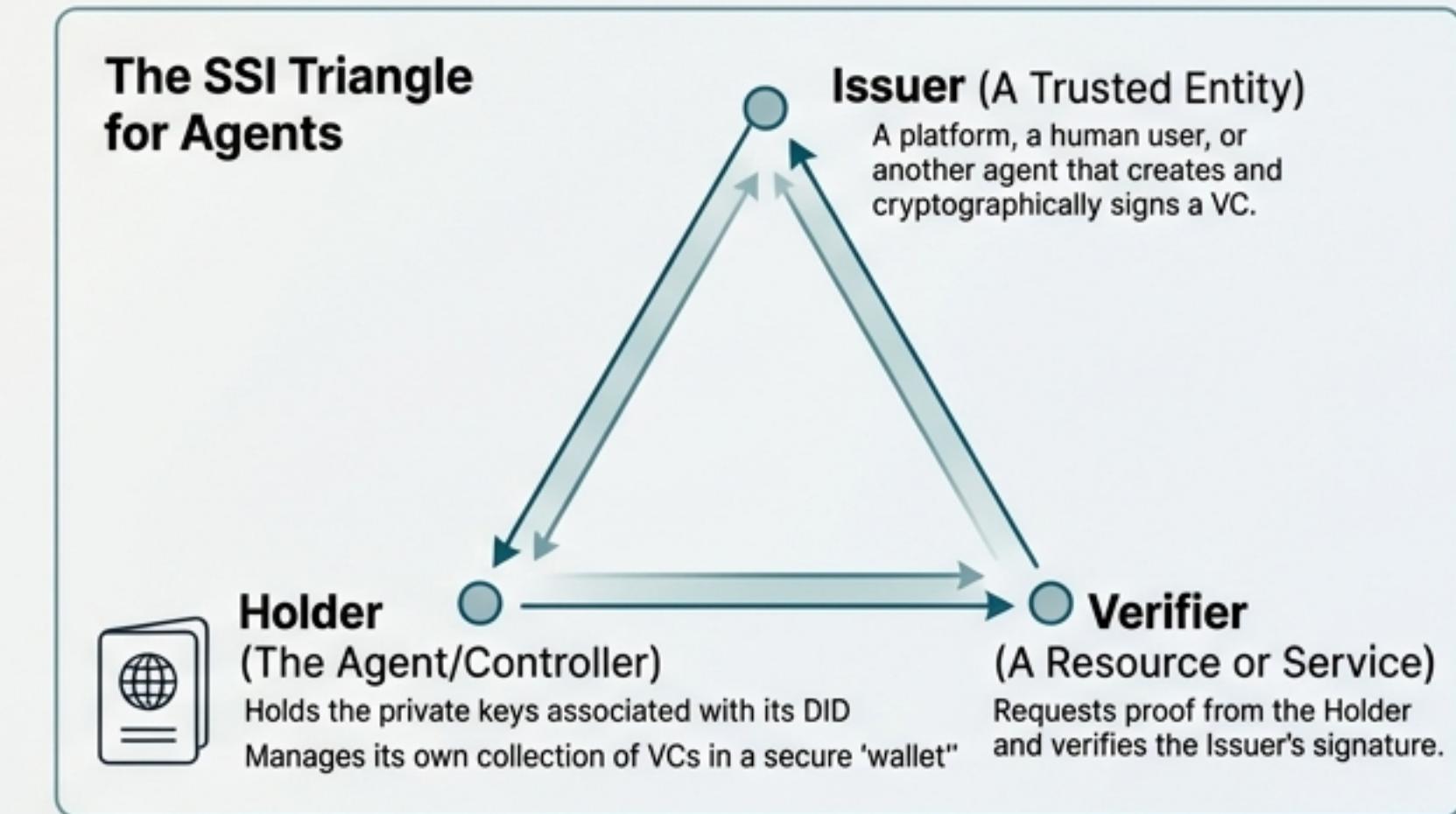
DIDs and VCs create a rich, machine-verifiable audit trail for every agent action. Instead of a log entry showing "service-account-prod executed API," you see "agent `did:...:alpha-001` executed API, authorized by VC `id:...:task-xyz` issued by `did:...:workflow-engine`." This level of granularity is essential for forensics and accountability.

A Paradigm Shift: Self-Sovereign Identity (SSI) for Agents

The agent's identity should belong to the agent (or its controller), not to a central provider. SSI applies this principle to create a portable, interoperable, and more secure identity model.



Old Model: Identity is 'rented' from a provider. If the provider goes down or revokes access, the identity is gone.



SSI Model: Identity is 'owned' by the agent. It is independent and portable across different systems and trust domains, just like a physical passport.



Why it Matters for Security Ops

SSI establishes a definitive, cryptographically certain line of accountability. When an incident occurs, the signed actions can be traced back to the specific agent that holds the corresponding private keys. This non-repudiation is critical for mitigating **ASI10: Rogue Agents**, as it ensures that every agent is uniquely identifiable and accountable for its behavior, preventing it from 'hiding in the crowd' of generic service accounts.

Dynamic Defense: Just-In-Time (JIT) Permissions with Ephemeral VCs

Enforce true least privilege by abandoning persistent roles. Instead, issue agents extremely narrow, time-bound Verifiable Credentials for a single task, which expire immediately after use.

The JIT Access Flow

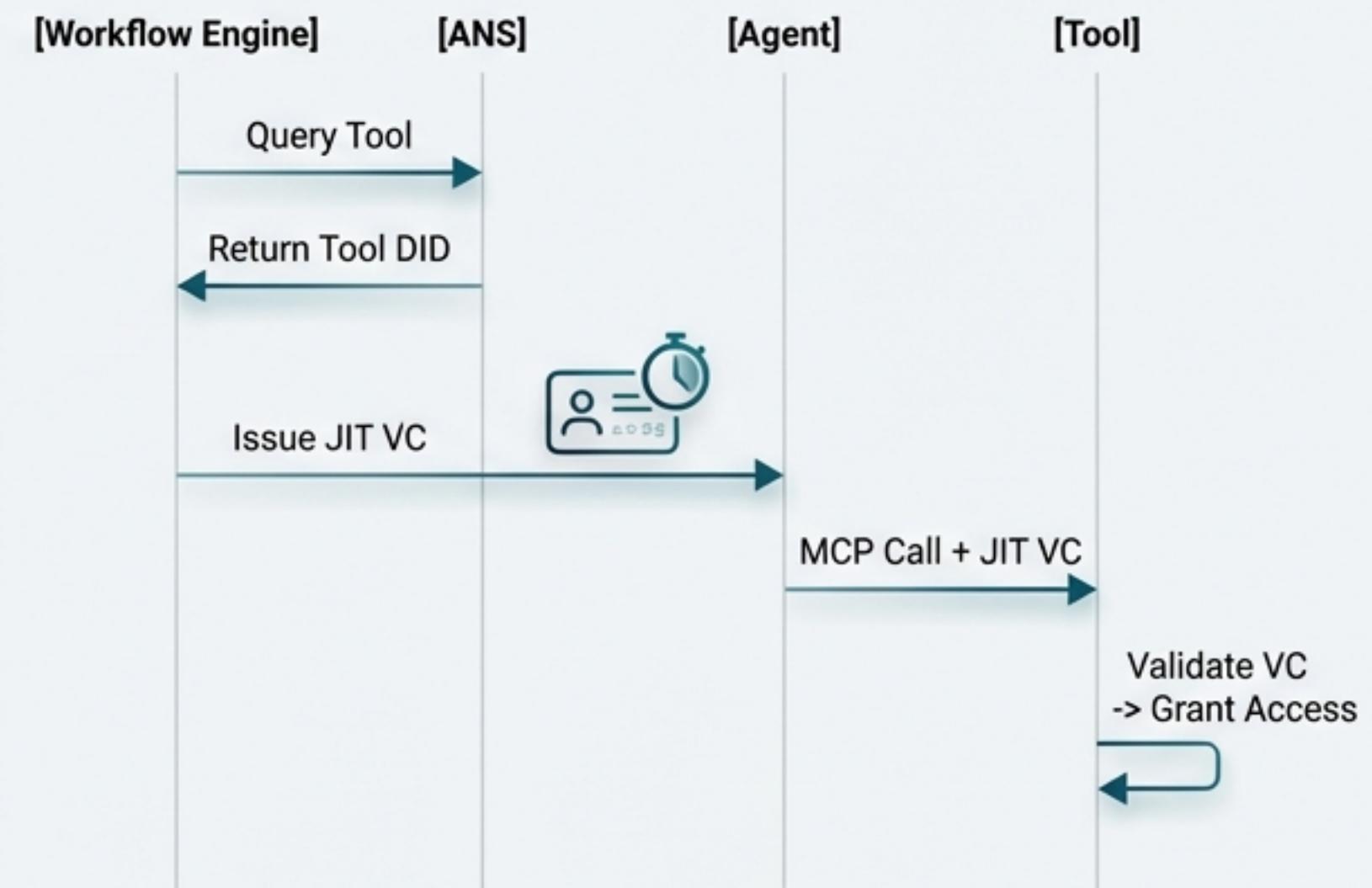
1. **Discovery:** A Workflow Engine queries the Agent Name Service (ANS) to find an available instance of DataTransformationTool-Q.
2. **JIT VC Issuance:** The Engine acts as an Issuer and generates a highly-scoped VC for the agent (did:ephemeral:task-xyz:agent-77).

```
{  
  "Subject": "did:ephemeral:task-xyz:agent-77",  
  "Authorized Tool": "did:com:acmetools...:instance03",  
  "Allowed Action": "executeTransform",  
  "Validity": "validUntil: 2825-10-02T14:45:00Z",  
  "Job ID": "job-ephemeral-77a"  
}
```



3. **Tool Invocation:** The agent presents this JIT VC within its call to the tool.
4. **Verification:** The tool's enforcement point verifies the VC's signature, validity period, and checks that the requested action matches the allowedActions and jobId. Access is granted only if all conditions are met.

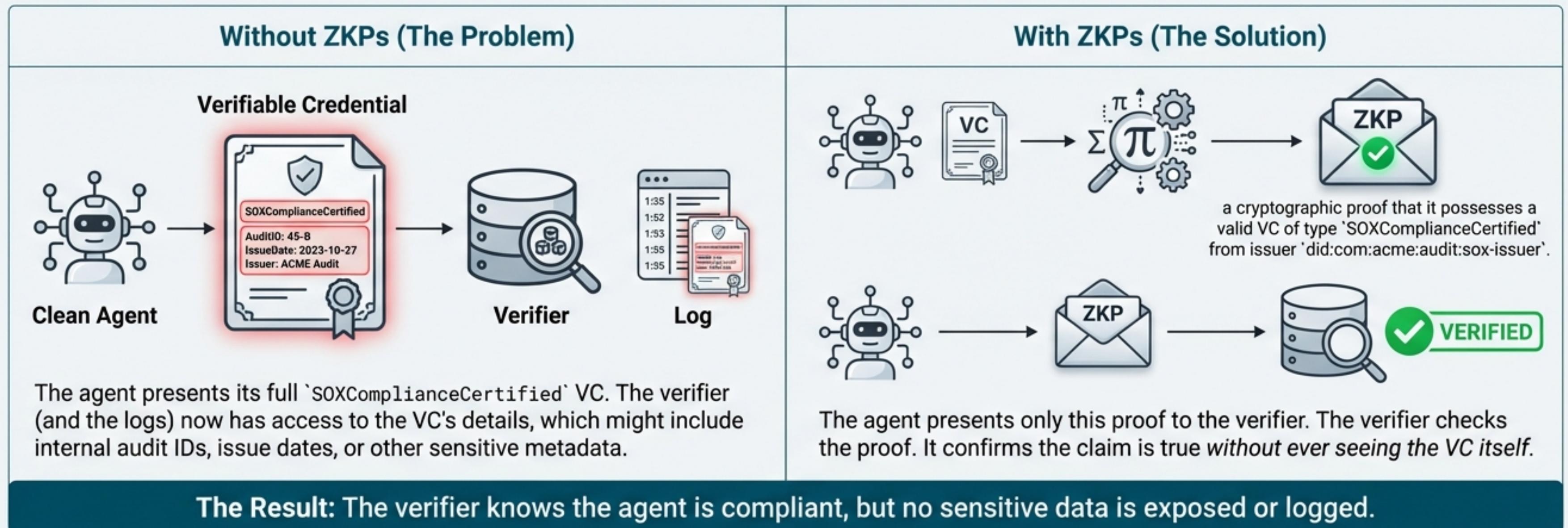
Visual Depiction



Why it Matters for Security Ops: JIT access dramatically shrinks the attack surface. A compromised agent or stolen credential is no longer a persistent threat. Its utility is limited to a few minutes and a single, specific action. This is the most effective defense against **ASI03: Identity & Privilege Abuse**, as it prevents both privilege retention and un-scoped inheritance by design.

Verifiable Compliance, Preserving Privacy: Zero-Knowledge Proofs (ZKPs)

Zero-Knowledge Proofs allow an agent to prove it possesses a required attribute or credential without revealing the underlying sensitive information, balancing verifiability with privacy.

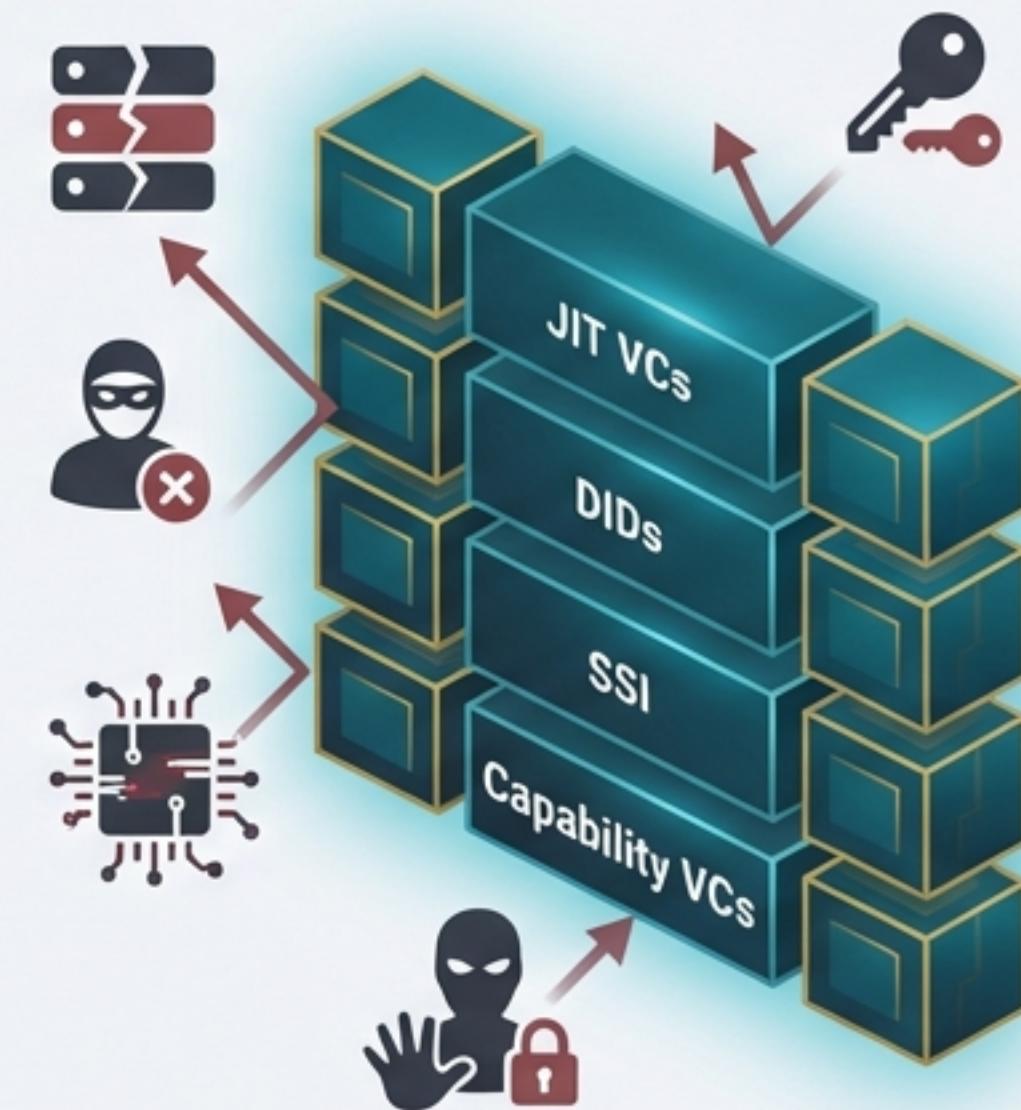


Why it Matters for Security Ops: In regulated industries (finance, healthcare), logging and auditing agent actions is non-negotiable, but exposing PII or confidential business information in those logs is a compliance nightmare. ZKPs enable a 'prove, don't show' model for policy enforcement. SecOps can build automated compliance checks that verify an agent's right to access data without creating a toxic repository of sensitive information, helping to meet mandates like the EU AI Act.

Building the Fortress: How Layer 1 Mitigates Top Agentic Threats

The components of the Identity & Credential Management layer are not just technical primitives; they are targeted defenses against the most critical threats facing agentic systems.

OWASP Agentic Threat	Layer 1 Mitigation	How It Works
ASI03: Identity & Privilege Abuse (Un-scoped inheritance, confused deputy, privilege retention)	JIT Verifiable Credentials	Issues task-scoped, time-bound tokens that expire after use, eliminating persistent privileges and shrinking the blast radius of a compromise.
ASI07: Insecure Inter-Agent Communication (Spoofing, Man-in-the-Middle attacks)	Decentralized Identifiers (DIDs)	Provides a unique, cryptographically verifiable identity for every agent. Mutual authentication using DIDs ensures agents are communicating with the intended party, not an imposter.
ASI10: Rogue Agents (Behavioral divergence, lack of accountability)	Self-Sovereign Identity (SSI)	Establishes a non-repudiable link between an agent's cryptographic keys and its actions. Every decision is signed, creating an immutable audit trail that makes accountability inescapable.
ASI02: Tool Misuse & Exploitation (Over-privileged tool access, unsafe delegation)	Capability-based VCs	Credentials explicitly define the exact tools and functions an agent is authorized to use (e.g., `toolName: SecureSQLConnector`, `targetSchemas: [Sales]`). This prevents an agent from using a legitimate tool for a malicious purpose.

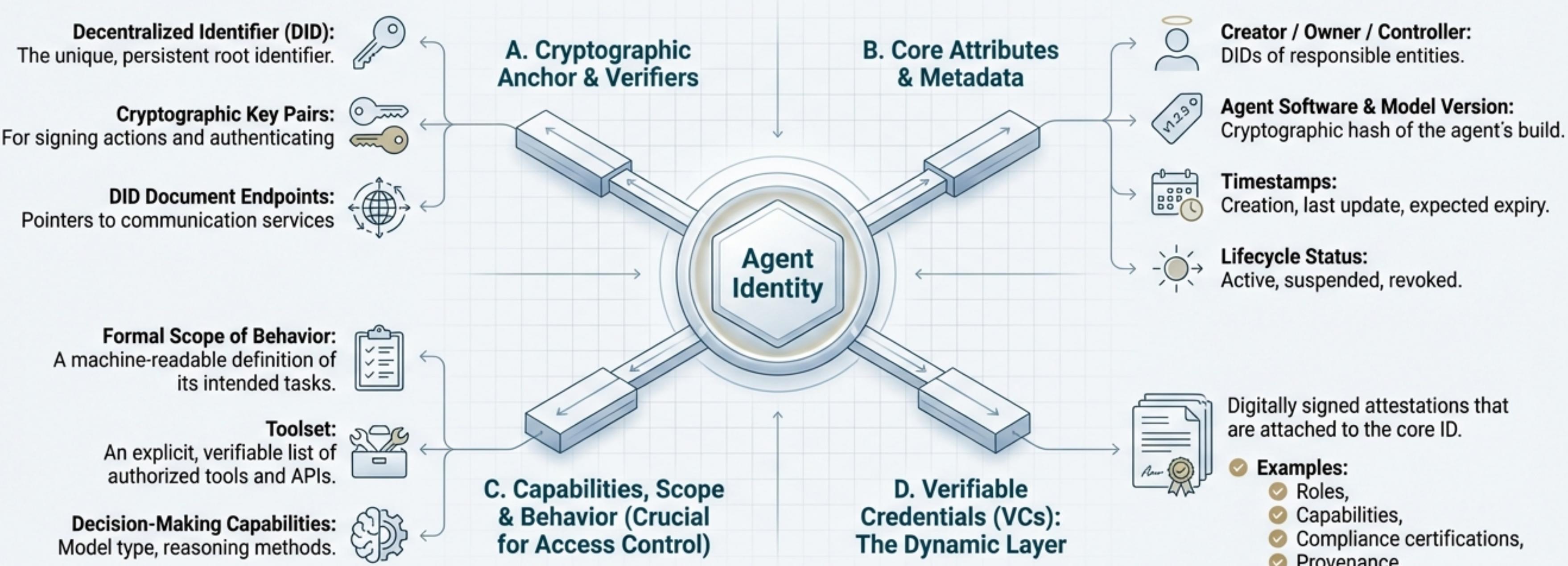


Why it Matters for Security Ops

This framework provides a cohesive strategy for moving from a reactive to a proactive security posture. By building security into the identity fabric itself, you can prevent entire classes of attacks by design, rather than trying to detect them after the fact.

The Anatomy of a Modern Agent Identity

An Agent ID is not a simple string. It is a rich, dynamic, and verifiable digital profile composed of multiple, interconnected components that provide a holistic view of the agent's purpose, provenance, and posture.

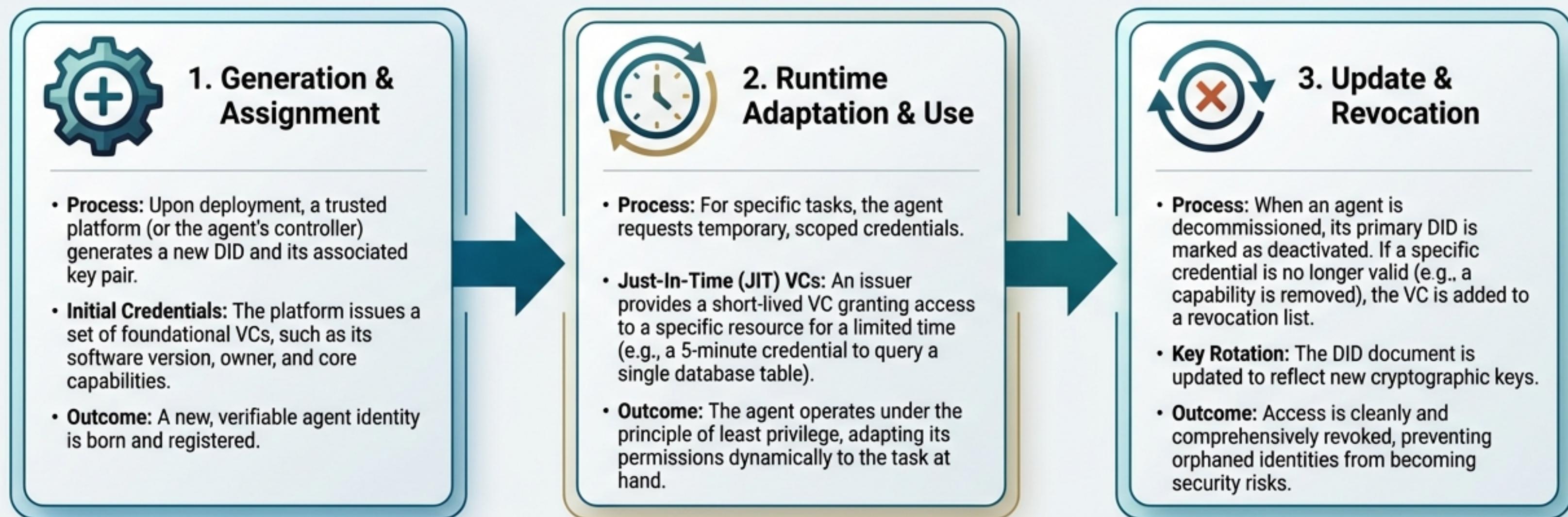


Why it Matters for Security Ops

This rich data model is the fuel for a sophisticated Policy Decision Point (PDP). Security policies can be based not just on a role, but on the agent's software version, its creator, the specific toolset it's authorized to use, and its verifiable compliance status. This enables truly fine-grained, context-aware access control.

Layer 1 in Action: Identity Lifecycle Management

A secure identity framework must manage the entire lifecycle of an agent's identity—from birth to revocation—to prevent orphaned credentials and outdated permissions.



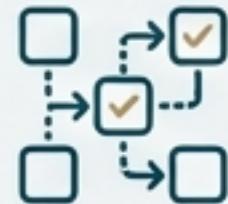
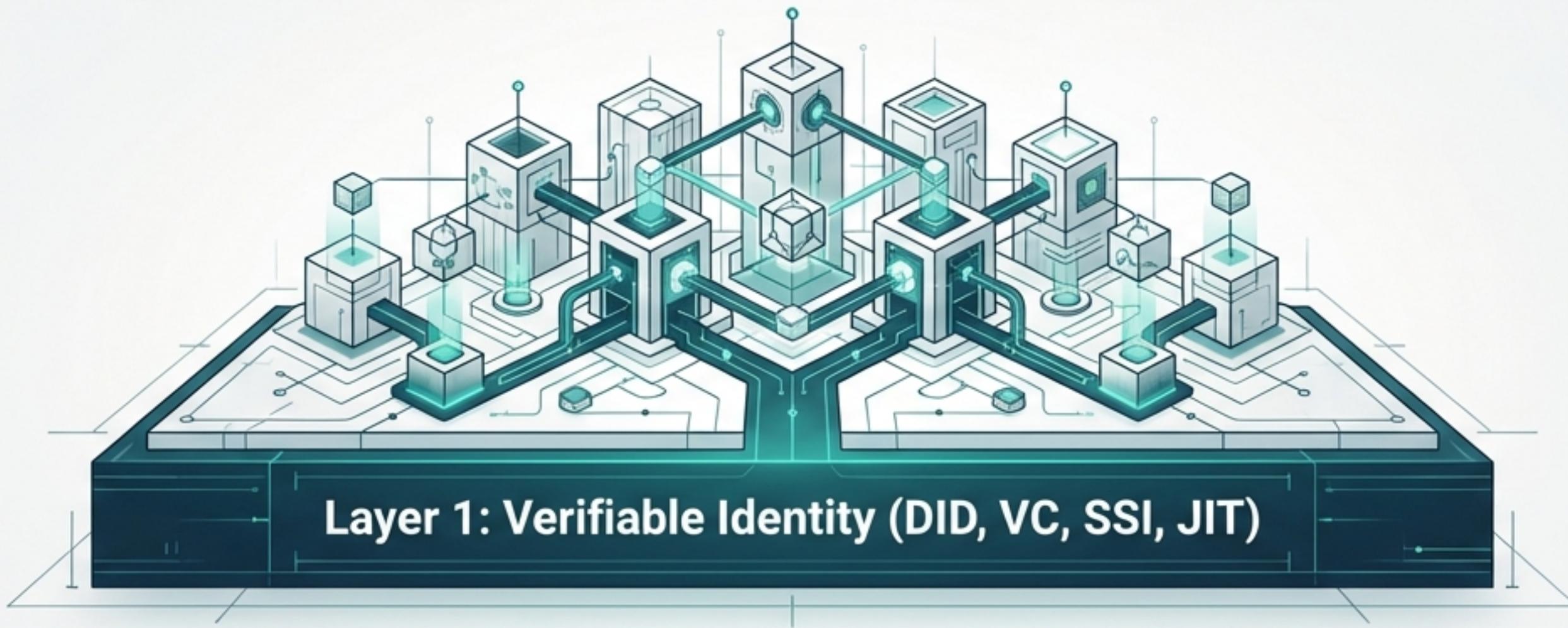
Why it Matters for Security Ops

Proper lifecycle management directly addresses **NHI1: Improper Offboarding**—a critical vulnerability where decommissioned identities retain access. This ensures that an agent's access rights are as ephemeral as the agent itself, eliminating a common source of persistent security holes.



"Traditional IAM systems provide a shaky foundation for the towering edifice of interconnected, autonomous AI agents."

A purpose-built, decentralized identity layer is not an optional add-on; it is the unavoidable prerequisite for building secure, accountable, and trustworthy agentic ecosystems.



From Centralized Roles to Decentralized Credentials

Move away from static, monolithic roles and toward dynamic, granular, and verifiable attestations of capability.



From Rented Identity to Owned Identity

Embrace Self-Sovereign Identity (SSI) to establish clear, non-repudiable accountability for every agent.



From Persistent Privilege to JIT Access

Make Just-In-Time permissions the default operational model to enforce true least privilege and minimize the impact of a compromise.