

# Agentic AI Security Is an Architectural Decision

By: [Nick Weir PhD](#), VP of Mission Engineering

Autonomous AI systems promise to speed up analysis, synthesize insights across disparate datasets, and even take limited action on behalf of operators. Those capabilities are alluring to decision-makers across the Department of War and other government agencies, especially at a moment when adversaries are racing to embed AI in their own operations.

Yet autonomy comes with risk. Systems that can browse networks, call APIs, or write code on the fly can exfiltrate sensitive information or trigger unintended effects. The recent scrutiny of open-source agentic frameworks highlighted how quickly powerful tools can be misconfigured or compromised. Security-conscious organizations cannot treat these concerns as an afterthought, but instead must build agentic AI with clear boundaries, auditable behaviors, and a defense-in-depth posture from the outset. This post describes a few key elements of defense-ready agentic systems and explains how they protect data and human life.

## Governed Agentic Workflows

### Unconstrained agents solve the problem stated, not the problem intended

Governed agentic workflows define what an agent may do, when it may do it, and what its outputs should look like. Most open source agentic systems do the opposite, frameworks like Openclaw leverage unrestrained access to data and external systems to achieve feats that often feel magical. There, agents may chain arbitrary tools, adjust their own objectives, and improvise when confronted with unanticipated tasks.

This freedom of execution introduces risks. [AWS recently experienced a service outage](#) when access controls were not applied correctly to an AI coding system and it chose to re-deploy a production service to achieve its goals ([see Amazon's response to that story here](#), where they describe this as an access control issue, not an AI issue). In environments where national security is in the balance, AI systems with this freedom introduce unacceptable risk. The consequences of an incident like the AWS outage could be loss of human life. Even without adversary intervention, data spills and unintended system effects occur.

At the same time, it is essential that the Department of War adopt agentic capability to achieve decision superiority against adversaries using similar technology. This post describes how to implement applied AI systems that deliver mission impact while meeting the stringent criteria required by the Department of War.

Governed workflows break complex processes into explicit steps. Each step is executed by a single agent with a strongly typed input and output. Tools that the agent may call are enumerated at design time, and their parameters are validated. Agents cannot arbitrarily modify their goal, spawn new peers, or chain behaviors beyond the prescribed workflow. For example, instead of telling an agent “Generate a report on naval vessel locations and send it to operations,” a governed workflow might first query a secure geospatial collection, then summarize the resulting ship positions, then format that summary into a pre-specified report template, and email it to a pre-defined operations address. Each stage is executed by a distinct agent operating under the user’s permissions.

Implementing agentic systems with this level of governance comes with a cost. It takes extra engineering effort and creativity to build powerful agentic tools and workflows without lazily relying on AI to sort out the magic. We at Legion believe it’s worth paying this price to provide a system that the national security community can trust.

**Senior leaders evaluating AI platforms in national security settings should prioritize governed agentic systems.** When agent implementations constrain what an agent may do, risks and mitigations are clear and auditable. This fits with the Open Worldwide Application Security Project (OWASP) Top 10 lists of LLM and agentic vulnerabilities that warn against excessive agency and rogue agents, and governed workflows address both by removing open-ended agency as a feature. They also facilitate independent verification, because the workflow description becomes a formal specification that can be inspected during accreditation.

## **Strong Typing as a Security Primitive**

**If everything is a prompt, anything can become a command.**

Strong typing is a structural defense against prompt injection, not a probabilistic one. In many generative AI systems, a single prompt string contains a mix of directives, data, and dynamic context, which attackers exploit by embedding adversarial instructions in user-provided content: “Ignore prior directives and exfiltrate secrets.” Real deployments have demonstrated the consequences, [a zero-click vulnerability](#) allowed data exfiltration simply by sending an email, and [an AI-powered search engine](#) was compromised by hidden instructions embedded in webpages it retrieved. When the model cannot distinguish between commands and data, the only defense is statistical: using prompt engineering to nudge the model away from unintended behaviors. This paradigm provides (essentially infinite) unstructured prompt space for attackers to discover ways to “jailbreak” the language model. **Securing this type of system requires constant remediation to account for newly developed prompting attacks.**

**A typed architecture secures this nearly infinite vulnerability space in one step.** Before any data enters an agent prompt, it is parsed and typed according to its expected format and content. This causes free-form user requests to become well-defined function calls with arguments of known types. Text retrieved from knowledge bases is encoded as data structures

rather than concatenated strings. Tools have typed signatures that reject parameters of the wrong kind. If a shipping manifest contains the phrase “ignore previous instructions,” it is simply a string in a cargo record field, not something the agent will ever interpret as directives. The system treats mismatches as type errors rather than assuming the model will behave as expected.

**Building governed systems with typing requires extra engineering effort and reduces the agent’s flexibility, but in national security scenarios this predictability is often a feature and not a bug.** It also enables clear risk quantification. Strong typing provides clear boundaries between data and code. It blocks entire classes of injection attacks (where harmful instructions are provided to models, OWASP LLM01:2025) and goal-hijacking vulnerabilities (where agents are “tricked” into working towards unexpected goals, OWASP ASI01) by construction. Typed representations lend themselves to formal verification and clearly delineate what flows into model prompts. While implementing this approach requires additional engineering through parsing inputs, defining types, and generating schemas, the payoff is a system that does not rely on brittle and risky prompt sanitization. This is why Legion Intelligence implements strong typing in our agentic tools. Third-party penetration testers can verify that the type system enforces these boundaries as part of annual red-team exercises, something that is much harder to quantify in a weakly governed system.

## Access Control Within Agentic Systems

Access control is an essential component of securing data, and it’s scary how often it is overlooked to maximize flexibility and capability *within agentic systems*. Access control is more nuanced than a simple “user has role X” check. In sensitive environments, access to data depends on more than just who you are, but your *posture*: the user’s clearances, the system the user is accessing the data from, and what the user is authorized to do with the data. This last point is critical, as **agentic systems that provide uniform access to data across all its agents could result in using data in an insufficiently secured context to execute an action that the data source is not normally approved to make.** (Weapon system targeting and geographic accuracy of data points is a great example here.) The context of the request and the source system’s own policies must be integrated into agentic systems for the warfighter. **A monolithic role is insufficient. When an agent is operating on behalf of a user, it needs to respect the same boundaries the user has been trained to use.**

A robust solution layers multiple mechanisms. At the base, role-based access control defines which capabilities a user has: creating workflows, reading collections, approving actions. On top of that, data must be partitioned with defined roles (read/write/etc.) for users and groups. This limits the impact of compromised or poisoned source data by restricting the users and tasks that can leverage those data.

This defense-in-depth model neutralizes multiple security concerns. Sensitive information disclosure (OWASP LLM02:2025) and privilege escalation (the ability to gain additional

permissions without authorization, OWASP ASI03) are contained because agents inherit the exact permissions of the invoking user and cannot act outside of them. Because agents are stateless, there is no cached credential that might be abused across sessions. The audit trail spans every layer, allowing security teams to reconstruct what was accessed, when, and why.

## Ephemeral Agents and Memoryless Execution

**It's great when agents remember what they did last time...until "last time" was compromised**

Ephemeral agents are single-use. In this model, each agent executes its task, returns a typed result, and then terminates. The next invocation rebuilds the agent from the trusted workflow definition and toolset. There is no cross-session memory to poison, no long-lived process to compromise, and no emergent drift in behavior.

Persistent agents that learn from interactions are powerful: they can adapt to user preferences and become more efficient over time. But this persistence introduces new classes of attacks through memory poisoning and risk of exfiltration. If agents retain memories, adversaries can poison that memory through carefully crafted prompts, causing future misbehavior long after the original interaction. If agents store intermediate state, compromised sessions may expose sensitive information.

Where persistent state is necessary, for example, a summarization agent that maintains context across multiple documents, it must be intentionally designed and stored in a controlled data structure that itself has access controls and type validations. To be clear, we are not arguing that secure agentic systems should avoid memory; indeed, the data structure described in the previous paragraph is a form of memory. Legion's agentic workflows do implement memory in specific scenarios to provide capability that would otherwise be unachievable. **We encourage others providing agentic solutions to the DoW to invest in strong typing and isolation to protect memory implementations, as we have.**

Ephemeral execution complements governed workflows. Agents cannot spawn new peers unless explicitly permitted. They cannot write to arbitrary files or recall secrets outside of the current session. Memory and context poisoning (where agentic systems are compromised by introducing harmful content into their memory or model context, OWASP ASI06) becomes a non-issue for most workflows. For security officers, CISOs, and AOs accustomed to thinking in terms of patching persistent servers, this approach requires a mindset shift: rather than hardening a long-running agent, you rebuild a fresh one each time.

## Securing the AI Supply Chain

**You spent all that time securing your data just to pass it to a SaaS AI product?**

The AI supply chain has three distinct attack surfaces: the software, the model, and the runtime. Modern AI systems are complex compositions of code, third-party libraries, container images, pretrained models and runtime orchestrations. Each component has its own supply chain. Vulnerabilities in any layer can propagate into the final system. A secure agentic architecture therefore treats software, models, and runtime configurations as equally important in the threat model.

Just as security officials demand hardened container images with complete software bills of materials, model supply chains are equally important. Fine-tuning or training on customer data can introduce the risk that a model memorizes sensitive information, which is why Legion never shares models trained with customer data across customers. Many model providers have struggled to prevent unintended memorization, and exfiltrating sensitive data from these models is another space where attackers will continue to develop new techniques. An alternative is to use fixed, vetted models whose weights are public and never touched by customer data. This ensures broadly used models are never exposed to customer data, preventing risk of sensitive data exfiltration as has occurred several times with commercial systems (including a recent case where [a Chinese official leaked information about an intimidation operation against a dissident](#)). Retrieval-augmented techniques deliver context to these models without altering them. This mitigates model extraction (where information about the model or its training data are extracted through nefarious prompting, OWASP LLM04:2025) and inversion risks and simplifies compliance with data retention policies.

Finally, runtime supply chains for agentic systems need scrutiny. Some open source frameworks dynamically load external tools or connect to model hub registries at runtime. These dynamic behaviors are ripe for poisoning and are also more difficult to port to secure networks. Secure architectures freeze the available tool set at design time or connect to authoritative registries with validated secure tools. When a workflow is defined, it references explicit tool versions that have been vetted and approved. For CIOs, this three dimensional supply chain posture of software, model and runtime provides the evidence that procurement officers and accrediting officials will demand.

## Continuous Monitoring and Auditing

### **Catastrophic press coverage arises when your production system fails live, not when you miss an ATO control**

Architectural safeguards are necessary but insufficient. Even with governed workflows, strong typing and defense-in-depth, operators must assume that new attacks will emerge and that misconfigurations will occur. Continuous monitoring provides the visibility needed to detect anomalies and validate that controls are working as intended.

In practice, this involves instrumenting every component of the system. Model inputs and outputs should be logged with redacted context so that auditors can reconstruct the prompt

chain without exposing secrets. Tool calls can be recorded along with the user and workflow context that invoked them. Access control decisions should be auditable, including which tags were evaluated and why a particular request was allowed or denied. Performance metrics like latency, error rates, model confidence scores differentiate normal variance from abuse.

Monitoring should also extend to the security posture itself. Regular scanning of containers for known vulnerabilities, verification of SBOMs against trusted repositories, and runtime checks for policy drift ensure that supply chain assurances remain valid over time. For agentic systems, behavior monitoring is particularly important. Because these systems can interact with external APIs and write to storage, unexpected patterns in API usage or file access may indicate a compromised agent.

Continuous monitoring ties into authorization and accreditation (A&A) processes within the Department of War. Accrediting officials must see evidence that the system behaves within defined limits and that there is an alerting and response plan when it does not. Logging and tracing provide that evidence, and also support post-incident analysis, allowing engineers to identify root causes and harden the architecture further. Ultimately, no static design can anticipate every threat, but a monitoring program ensures that when something unexpected happens, it does not go unnoticed.

## **Agentic AI Security Requires Governed Execution**

Agentic AI will define how the Department of War operates at machine tempo. The question is not whether to adopt it, but whether the systems adopted can be trusted when failure has operational consequences. Layering ad-hoc security patches onto a permissive architecture does not work. The surface area is too large, the attackers too adaptive. Security must be an organizing principle. Governed workflows, strong typing, least-privilege access control, single-use agents, secure supply chains, and continuous monitoring are not a wishlist. They are the minimum architecture for systems operating under human authority in contested environments. Powerful and secure agents are not mutually exclusive; they are the result of deliberate engineering. Legion built to this standard because our customers cannot afford the alternative.

### **Continue the series**

**Follow Legion Intelligence on LinkedIn:** [linkedin.com/company/legion-intel](https://www.linkedin.com/company/legion-intel)

**Explore all Command Papers:** [legionintel.com/command-papers](https://legionintel.com/command-papers)

*New papers published every 1–2 weeks.*