

# Trustworthy AI Is Not a Feature. It Is the Prerequisite.

A Governance Framework for Enterprise Agentic AI

---

Clay House · Cognita AI Inc. · March 2026

In order for us to drive our cars, fly in an airplane, cross a bridge, take medicines, or eat the food we consume, we must trust the process of delivering those products and services to us. Without that trust, we would do everything in our power to find alternatives.

We have a baseline level of trust in these things - not because the companies or the people involved are inherently trustworthy – but because the results are controlled to produce appropriately consistent and reliable outcomes within a risk/impact framework.

Building this trust requires:

- › **Responsible** decisions regarding the delivery of services and/or products with governance and oversight.
- › **Protection** of the services/products from nefarious actions such as tampering, compromise, and misuse.
- › **Controlled** execution preventing unintended impacts beyond the defined scope of authority and intent.
- › **Monitoring** of results against expected outcomes through audits, QA processes, and automated measures.
- › **Mitigation/remediation** of the root cause and effect of unintended outcomes.

A critical barrier to the adoption of AI solutions is achieving this trust. Not because of the unique challenges presented by AI, but because the foundational controls for building trust have significant gaps.

*We do not trust these industries because of confidence in their technology. We trust them because their results are controlled. Trustworthy AI requires the same foundation.*

To achieve Trustworthy AI, we must address all of the five components. A weakness in any layer compromises or destroys the overarching trust.

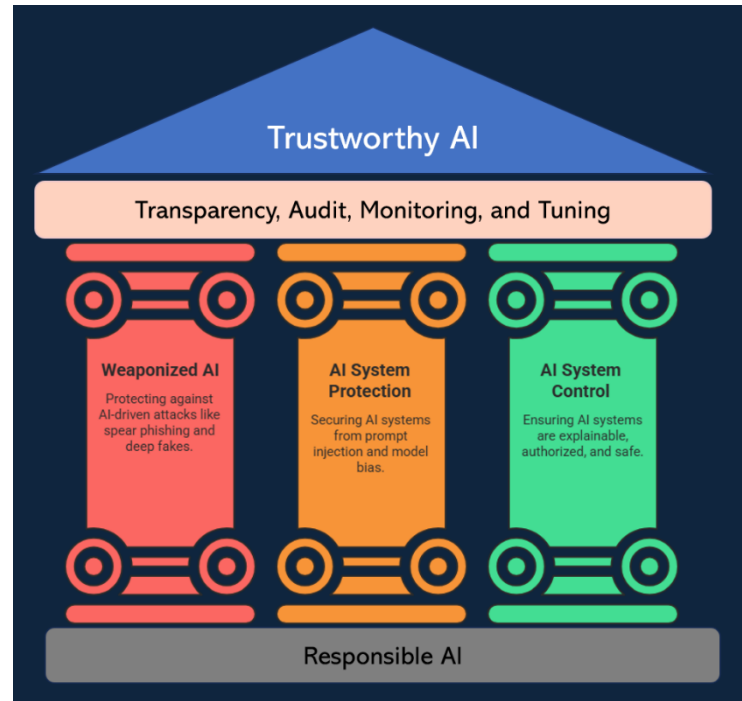
The foundation of trust is demonstration and governance over the responsible use of AI. Much like the use of clinical data must go before an Internal Review Board, clinical policies have Clinical Review Boards, and cars have engineering, crash, and regulatory reviews, we must implement the equivalent for AI

The responsibility foundation is owned by the business and informed by IT. The business must assess the risk, likelihood, and impact of incorrect AI outcomes to determine if its use is appropriate. If so, it must have continual oversight informed by the system's performance, changes in stakeholder expectations, and regulatory shifts to continually validate this decision.

The control columns address the ongoing operation of the system. Cybersecurity retains its traditional role of ensuring the Confidentiality, Integrity, and Availability of the systems.

Integrity takes on new meaning: the traditional protecting from external compromise of integrity expands to include business accountable functions such as protecting against inappropriate bias, adequate training data, system prompts, model tuning, etc. to ensure that the system continues to operate as intended. The business and IT must collaborate on securing the training pipeline and sources against poisoning, instruction manipulation, and other emerging threats.

Similarly business and IT must collaborate on controlling the behavior of the AI system. This is particularly important with Agentic AI as those solutions autonomously plan their execution, discover available resources, take actions, and evaluate to determine completeness.



A unique characteristic of AI systems is the black-box nature of their execution. Continuous monitoring of AI system flows, decision points, data used/excluded in decisions, and varying techniques from audits to statistical analysis are required to drive visibility into AI systems. It is the business that has the expertise and knowledge to interpret the collected data with IT being an enabler.

---

*It is impossible to trust a system when you have no ability to assess the validity of its results.*

---

decision points, data used/excluded in decisions, and varying techniques from audits to statistical analysis are required to drive visibility into AI systems. It is the business that has the expertise and knowledge to interpret the collected data with IT being an enabler.

## AI Systems

In traditional systems, the logic for the system is expressed using a programming language. This makes the system behavior knowable (e.g. deterministic) by examining the programming instructions. This is great for pre-defined tasks but cannot handle general purpose tasks that can't be readily programmed.

Instead of being given a list of rules, AI systems are **trained**. To learn how to program, an AI analyzes vast amounts of existing code—like a student reading every textbook and project ever written.

The AI turns all that information into a complex mathematical map. When you ask it a question, it doesn't 'think' in the human sense; it looks for the mathematical patterns in your request and finds the most likely 'next piece' of the puzzle based on its training.

It predicts the answer by matching these patterns. Because the system is based on **probability** rather than rigid rules, it might give slightly different answers to the same question—much like how a person might explain the same concept using different words each time.

In simplistic terms, you can think of AI systems as falling into one of two categories:

- AI Enabled
- Agent-based (Agentic AI)

An AI enabled system uses AI to perform tasks that are not well suited to traditional programming such as pattern recognition (images, text, programming, etc.), summarization of material, chat bots, language translation, etc.

Agentic AI systems are ones where we delegate the execution of a task, and its subtasks, to an AI enabled “agent”. The agent uses its knowledge of the assigned task and any given constraints. It develops an execution plan, identifies required resources such as other agents, tools, or data, executes the plan, and evaluates the result to determine if the task is complete.

## Why Agentic AI Makes This Harder — and More Urgent

Agentic AI presents a number of new challenges:

- We can provide “hints” as to steps but nothing ensures the AI follows all steps or in the same manner.
- The AI must interpret given and retrieved data consistently with the task. Similar sounding data names, data with multiple dates but no clear description, etc. increase the challenge.
- AI can get “distracted” by irrelevant data
- Tool selection and use is not guaranteed to occur or occur properly
- The agents, resources, data, etc. accessed are not pre-defined and cannot be attributed

Traditional controls and governance models were not designed for this behavior. For example, it is not sufficient or even possible to define what a tool can/cannot do if you have no idea who is using it or why. That appropriateness must be derived at runtime based on the context of the action.

This is what makes AI such a powerful a tool. We use AI enabled systems to perform tasks beyond traditional programming and Agentic AI to act as workers for broader tasks that require steps and reasoning.

The **control objectives** of our traditional control frameworks still hold in concept but the meaning, importance, and how they are achieved materially change. AI also adds additional responsibilities due to the probabilistic nature that must be owned by the business to detect, understand, and address behavioral deviations that are inherent in the technology, not rooted in malicious behavior or compromise. Integrity now extends well beyond its traditional scope in Cybersecurity to reflect the integrity of the business intent of the system ensuring that it is operating faithfully within the business authority and intent. In this context, business authority is not controlling access in context but in ensuring the system doesn't exceed the business authority. This is equivalent of an employee weighing in on a discussion or decision for which they are not authorized nor qualified.

### The Deterministic Assumption

Traditional controls such as defining access rights, testing, quality checks on results, etc. are designed with the expectation that they are always used in exactly the same way. Because of this we can reasonably enumerate the various scenarios, test for them, design access control for them, and audit against expectations.

AI agents violate this assumption structurally. They do not execute predefined logic. They infer actions from context. Given identical inputs, an agent may select different tools, call different services, interpret data differently, and construct a different work plan. Even when outcomes are consistent, the reasoning path that produced them is opaque and non-reproducible. The same agent, asked the same question tomorrow, may take a meaningfully different path to the answer.

As an example, think of a simple utility (e.g. tool) that can read, write, modify, and delete files.

One can easily both see the necessity of such a tool as well as how it can be extremely dangerous. The reality is that you can't predetermine what this tool should be allowed to do or not to do. Instead you must evaluate it in the context of its use – who authorized the task, what is the business intent of the task, and what are the specific details of the task.

*In AI systems, **control** and **integrity** refer not only to what the system can do, but what it can “**say**”.*

Consider a developer using AI agents for code development. When writing or editing code, this tool would require read/write/edit access. To determine what project files they can change depends on the authorizer of the task (the developer). Most AI coding tools have “planning” vs “editing” modes where in the former access is read only. The context of

- Authorizing Authority – A specific developer determining which projects/files can be accessed
- Task – Planning vs Editing
- Context – The specific objectives defined by the developer, coding standards defined by the enterprise, active project, etc. shape the behavior, guardrails, and output

## Control Aggregation: Invariant + Task + Process + Context Bounded by Scope of Authority

### Ground Truths/Characteristics

Agents, tools, and resources very likely have a ground truth of controls and characteristics that should be invariant – never change regardless of context. A crude, but obvious example, is that basic safety rules of power tools apply regardless as how the tools are used. As with all controls, there will be a range of risk, likelihood, and impact. Those controls will always apply but may be modified by the usage and context (e.g. using an electric power tool in dry open conditions vs wet vs in the presence of explosive gasses).

### Task

These are core governing rules and controls that apply within a specific task. Continuing the tool analogy, there are different safety considerations when using a hammer to drive a nail vs extracting a nail. Task controls are invariant to the task regardless of why the task is being performed - e.g. building a house or putting up shelves.

### Process

A process is a collection of tasks in the context of a specific goal. Shifting gears to a more business-focused example, consider a healthcare eligibility check. There are base invariants like the definition of a patient, verification of insurance coverage, etc. However there are many process-specific constraints:

Process	Characteristics	Constraints/Controls
HIPAA Eligibility	Only confirms insurance coverage and basic benefit structure.	Must adhere to Federal HIPAA patient matching rules and information returned.
Prior Authorization	Eligibility is extended to look at detailed benefit design, accumulators, medical necessity, etc.	Defined by healthplan medical policies and benefit designs within state and federal regulations
Claim Acceptance	Confirms insurance coverage for date of service and Coordination of Benefits with claim structure validation	NOT bound to HIPAA patient matching rules allowing for soft matches and scenarios such as attributing a newborn to a mother's coverage.

There is no defined set of access rights, control checks, etc. for an “Eligibility” agent absent the context of the task within a business process. These could certainly be built as different eligibility agents/services with the correct controls defined within those implementations which is a valid approach. However the gap remains that an autonomous agent would need to definitively understand which agent/service to use in which context. Today that's not possible.

#### THE CONTROL GAP

*Traditional controls attach to execution paths. Agentic systems do not have fixed execution paths. The control frameworks organizations rely on – role-based access, audit logs, change*

*management, policy enforcement — were not designed for systems whose behavior cannot be predicted from their inputs. Adapting those frameworks is necessary but not sufficient. New infrastructure is required.*

## Introducing SADAR

The control gap is not a gap in organizational will or policy intent. It is a gap in infrastructure. Enterprises cannot govern what they cannot identify, and they cannot attribute what they cannot trace. Before authorization can be enforced, before audit trails can be constructed, before a runtime system can evaluate whether an agent action falls within its scope of authority — there must be a semantic foundation that makes those determinations possible.

SADAR — the Semantic Agent Discovery and Authorization Registry — is an open standard designed to provide that foundation.

SADAR defines how agents, tools, and resources are described, discovered, and attributed within enterprise ecosystems. It establishes a machine-readable registry model grounded in existing industry standards, enabling agents to identify capabilities unambiguously, carry verifiable identity into every interaction, and operate within governance boundaries that are defined, enforceable, and auditable.

Critically, SADAR does not replace existing governance frameworks. It provides the semantic infrastructure those frameworks require to function in agentic environments. The five components of Trustworthy AI remain the objective. SADAR makes them achievable.

The standard is published and maintained by OpenSemantics.org as an open specification, available for implementation across platforms and vendors.

**Trustworthy AI is not a feature. It is the prerequisite. SADAR is the infrastructure that makes it possible.**

## We Already Require These Controls — Over Humans

The most clarifying question in AI governance is not 'is AI trustworthy enough to deploy?' It is: 'what controls do we already require of humans executing the same processes, and why would we accept less from an AI system doing the same work?'

Consider how a mature healthcare operation manages its human workforce. Every call center agent operates against a defined process — scripted steps, required disclosures, escalation triggers that fire at specific conditions. Calls are recorded. QA teams sample those recordings and drop into live calls. Screen activity is captured. The data presented to the agent during the interaction is reviewed. Escalation paths are defined and enforced: agent to lead to senior lead to manager to director to VP to compliance or medical review, depending on the nature and severity of the issue. None of this exists because the workforce is incompetent. It exists because the process requires it — because the organization has an obligation to the people it serves, and because 'our people are well-trained' is not a substitute for demonstrable control.

These are not exceptional requirements for a high-risk industry. They are the standard operating model for any organization that takes its obligations seriously. Financial services has the same structure. Insurance has it. Any regulated industry with consequential customer interactions has some version of it.

Now consider what most organizations do when they deploy an AI agent to execute the same processes. The agent operates without a defined escalation path. Its decisions are not sampled. Its inputs and outputs are not reviewed against a process standard. There is no QA function evaluating whether it is following the defined procedure or drifting from it. The data it acts on is not audited for semantic consistency. There is no mechanism to detect when it should have escalated but did not.

This is not a gap that better models will close. A more capable agent making more confident decisions in a governance vacuum is a larger risk, not a smaller one. The question is not whether the model is good enough. The question is whether the organization has built the equivalent of the call recording infrastructure, the QA sampling function, the escalation trigger definitions, and the process audit capability — for its AI workforce.

*We would not accept 'our staff are well-trained' as a substitute for call recording, QA sampling, and defined escalation paths in a human contact center. Why would we accept it for an AI agent doing the same work?*

### The Observe-Plan-Act-Evaluate Loop and Its Governance Obligations

Each phase of the agentic loop creates distinct governance obligations that existing frameworks do not address.

Phase	What the Agent Does	Governance Obligation
Observe	Reads data from systems, interprets context, builds a model of the current state.	
Plan	Selects tools, sequences steps, determines how to accomplish the objective.	
Act	Executes calls against real systems, modifies state, invokes other agents.	
Evaluate	Assesses outcomes, decides whether to continue, revise, or escalate.	

### Why Stakeholders Cannot Be Asked to Simply Trust

A governance framework for trustworthy AI is not primarily a technical question. It is an accountability question. The organizations deploying AI agents bear responsibility for their behavior — to their customers, to regulators, to their own boards and employees, and to the partners and counterparties whose systems those agents interact with.

The current practice in most enterprises is to deploy agentic AI in constrained contexts, monitor it closely, and rely on human review to catch problems. This approach is honest about the risk, but it trades away the autonomy that justifies the investment. An agent that requires human review of every consequential decision is not an autonomous agent. It is an expensive suggestion engine.

The alternative — granting meaningful autonomy without the governance infrastructure to support it — creates a different problem. When something goes wrong, and in probabilistic systems operating at scale, something eventually will, the question is not just what happened but whether the organization took reasonable precautions. Deploying autonomous AI in consequential contexts without a demonstrable governance framework is not a defensible position before a regulator, a court, or a customer who has been harmed.

*The organizations that move fastest on agentic AI without governance infrastructure are not gaining a competitive advantage. They are accumulating liability.*

The table below captures what the principal stakeholder groups need from an AI governance framework — not as a wishlist, but as the minimum basis for reasonable acceptance of AI-influenced decisions in their domain.

Stakeholder	What They Need from a Trustworthy AI System
Customers	Informed consent for AI-influenced decisions. Recourse when AI produces harmful outcomes. Confidence that their data is used only within the scope they agreed to.
Regulators	Demonstrable governance process. Audit trail sufficient to reconstruct any decision. Evidence that controls were in place before deployment, not retrofitted after an incident.
Board & Executives	Risk exposure that is characterized, bounded, and within stated risk appetite. Liability that is defensible — not just 'we relied on the vendor.'
Employees	Clarity on where AI operates autonomously versus where human judgment is required. Protection from accountability for outcomes they had no ability to oversee.
Partners & Counterparties	Confidence that AI operating on their behalf or interacting with their systems is operating within agreed boundaries. Attribution when something goes wrong.

## The Governance Model We Already Have: The IRB Analogy

The most useful existing analogy for AI governance is not in technology. It is in clinical research. Institutional Review Boards exist because the potential benefits of clinical research do not automatically justify the risks to human subjects. An IRB does not prevent research from happening. It establishes a process by which proposed research is reviewed against defined criteria — risk to subjects, necessity of the methodology, adequacy of disclosure, capacity for meaningful consent — before it proceeds. It requires a decision record. It has authority to stop or modify research that does not meet the standard. And it operates continuously, not just at the point of initial approval.

The analogy to AI governance is close enough to be instructive. The potential benefits of agentic AI do not automatically justify the risks to the people and organizations it affects. An AI governance board modeled on the IRB does not prevent AI deployment. It establishes a process by which proposed deployments are reviewed against defined criteria — risk to affected parties, adequacy of the governance controls, scope of autonomous authority, disclosure requirements — before they go into production. It requires a decision record that is defensible after the fact. And it operates continuously, because AI systems in production evolve in ways that require ongoing oversight.

**A CLINICAL PARALLEL**  
*A hospital would not deploy an autonomous clinical decision system — one that observes patient data, forms a diagnostic plan, and initiates treatment without physician review — without demonstrating to a clinical review board that the system's decisions are grounded in evidence, its behavior is auditable, its scope is defined, and patients are appropriately informed. The bar for autonomous AI in consequential enterprise contexts should be no lower.*

The key features of the IRB model that are worth preserving in an AI governance analog:

<b>Prospective review</b>	The governance review happens before deployment, not after an incident. The organization demonstrates adequate controls at the point of approval, not at the point of defense.
<b>Defined criteria</b>	Acceptance is based on specific, documented criteria applied consistently — not judgment calls made ad hoc. Risk appetite, disclosure requirements, and control adequacy are defined in advance.
<b>Decision record</b>	Every material AI deployment has a documented governance review with a recorded decision. This record is the evidence of reasonable precaution.
<b>Continuing review</b>	Approval is not permanent. Systems that change materially, or that operate in contexts that change, are subject to re-review. Monitoring feeds back into governance.
<b>Authority to stop</b>	The governance body has actual authority to halt or modify deployments that do not meet the standard. A governance function without enforcement authority is an advisory function.

## Disclosure

Affected parties are informed of AI involvement in decisions that affect them, in a manner appropriate to the context and the stakes.

## A Risk-Based Framework for Trustworthy AI

Trustworthy AI is not a single control or a single product. It is a layered framework in which each layer addresses a distinct aspect of the accountability problem. The five layers below map to each other and to the existing governance frameworks that enterprises already operate. They are designed to be risk-based: the depth of implementation at each layer is proportional to the autonomy of the system, the sensitivity of the context, and the stakes for affected stakeholders.

<b>01</b> <b>Responsible Use</b> <i>The governance body</i>	An IRB-equivalent review process that defines acceptable AI use relative to the organization's risk appetite. Reviews proposed deployments against defined criteria. Requires disclosure to affected stakeholders. Produces a decision record for each material deployment. Has authority to stop or modify a deployment.	<b>ANALOGOUS CONTROLS</b> <i>Institutional Review Board (clinical research) · Clinical Policy Review Board</i>
<b>02</b> <b>Protect</b> <i>Integrity of the system</i>	Ensures that the components operating in the AI system are what they claim to be and have not been tampered with. Signed manifests, verified identities, and semantic data binding guarantee that the system executing is the system that was reviewed and approved.	<b>ANALOGOUS CONTROLS</b> <i>Supply chain integrity · Code signing · Change management controls</i>
<b>03</b> <b>Control</b> <i>Permissions envelope</i>	Defines what the system is permitted to do, with whom, and under what constraints. Discovery scoped to trusted providers. Bilateral policy enforcement. Compliance attestations enforced before connection. Rate limits, cost constraints, and data sovereignty requirements enforced at the infrastructure layer.	<b>ANALOGOUS CONTROLS</b> <i>Least privilege · Role-based access control · Vendor risk management</i>
<b>04</b> <b>Monitor</b> <i>Continuous visibility</i>	Transaction lifecycle capture — selection, inputs, and outputs recorded against the process definition. Feeds statistical QA against defined process standards, using the same methodology applied to human processes in call centers, clinical audit, and financial sampling. Makes remediation, tuning, and optimization data-driven.	<b>ANALOGOUS CONTROLS</b> <i>Statistical process control · Clinical audit · Model risk management</i>
<b>05</b> <b>Explainable</b> <i>Reconstructible decisions</i>	The ability to reconstruct and justify any agent decision in terms a human reviewer, auditor, or regulator can evaluate — not by opening the model's black box, but by answering: what process governed this, what did each participant claim to do, what was selected and why, what did the data mean, and where in the process did this occur.	<b>ANALOGOUS CONTROLS</b> <i>Audit trail · Model explainability · Right to explanation (GDPR Art. 22)</i>

Two aspects of this framework are worth emphasizing because they are frequently treated as afterthoughts in AI governance discussions.

## **Semantic and Process Integrity Are Security Properties**

Classical security frameworks focus on confidentiality, integrity, and availability — where integrity means data has not been tampered with in transit or at rest. In an agentic AI context, this definition of integrity is necessary but not sufficient.

An agent that reads a data field without tampering, but misinterprets it because the field means different things in different systems, has violated integrity in the sense that matters for business outcomes. An agent that executes all steps of a workflow with valid credentials and signed components, but skips a validation gate because no infrastructure enforced it, has violated process integrity. Neither violation shows up in a classical security audit. Both can cause material harm, and both are attributable to the organization that deployed the system.

Semantic integrity — the guarantee that data arguments mean what participants in a transaction agree they mean — and process integrity — the guarantee that defined business processes are enforced at the infrastructure layer rather than assumed to be implemented correctly in each component — are therefore first-class governance requirements, not enrichment features.

## **Explainability Is an Architectural Property, Not a Model Property**

The most common misunderstanding in AI explainability discussions is the assumption that explainability requires opening the model's black box — understanding why a neural network produced a particular output. This is technically difficult and in many architectures currently impossible.

But regulatory and governance explainability does not require this. It requires the ability to reconstruct and justify any system decision in terms that a reviewer, auditor, or affected party can evaluate. For an agentic system with appropriate governance infrastructure, this is achievable without ever inspecting model internals: what process definition governed this transaction, what did each participating component claim to do, what was actually selected and why, what did the data mean in the context of the governing standard, and where in the defined process lifecycle did this decision occur.

This form of explainability is an architectural property of the governance infrastructure, not a property of any individual model. It can be achieved today, in production systems, with existing technology — provided the governance infrastructure is in place to capture and expose the required record.

*You do not need to open the model's black box to explain an agent's decision. You need a governance architecture that captures what the agent was supposed to do, what it actually did, and what the data it acted on meant.*

## Challenges with the Current State

Most enterprises deploying agentic AI today are operating without most of this framework in place. They have models. They have frameworks for orchestrating agents. They may have logging. What they typically lack is: a formal governance review process for AI deployments, verifiable identity and integrity assurance for the components operating in their agent systems, semantic grounding for the data those agents consume, process integrity enforcement, and a cross-framework audit trail sufficient to reconstruct decisions after the fact.

This is not a criticism of the organizations involved. The governance infrastructure for agentic AI does not yet exist as a mature, standardized, deployable solution. Organizations are making reasonable decisions under genuine uncertainty. The risk is that 'we deployed before the standards existed' becomes a less defensible position as the standards develop, the deployments scale, and the first significant incidents surface.

The choice organizations face is not between deploying AI with governance and deploying AI without it. It is between three positions:

<b>Constrained deployment</b>	Limit autonomous AI to contexts where the governance gap does not create material risk — tightly scoped, human-reviewed, low-stakes. Responsible, but forfeits the business value that justifies the investment.
<b>Ungoverned deployment</b>	Deploy at scale without the governance infrastructure. Captures near-term value. Creates accumulating liability and regulatory exposure that is not visible until an incident or an audit makes it so.
<b>Governed deployment</b>	Build the governance infrastructure in parallel with capability deployment. Requires investment and discipline. Is the only position that is defensible to all stakeholder groups over time.

## What Governed Deployment Requires Technically

The governance framework described in this paper is not a policy document. Policies without technical enforcement are aspirations. The five-layer framework requires technical infrastructure at layers two through five — and that infrastructure does not currently exist as a mature, standardized component of the enterprise AI stack.

The Protect layer requires verifiable identity for every component in the agent system, signed manifests that attest to what each component is and what it is authorized to do, and semantic grounding that binds the data agents consume to published, external standards.

The Control layer requires that policy constraints — who can discover what, under what compliance attestation, subject to what rate limits and cost constraints — are enforced at the infrastructure layer before connections are established, not dependent on each component implementing policy correctly.

The Monitor layer requires a transaction record that is framework-agnostic and organization-spanning — one that survives across the boundaries between different agent frameworks, different vendors, and different organizations participating in the same workflow.

The Explainable layer requires that the process definitions governing agent behavior are first-class artifacts that exist independently of any agent implementation — authoritative blueprints that agents execute against, and that auditors can reference to evaluate whether execution was consistent with intent.

#### **THE INFRASTRUCTURE GAP**

*The governance framework described here is achievable with current technology. What is missing is not the technical capability but the standardized infrastructure layer that implements it in a vendor-neutral, interoperable way. The Semantic Agent Discovery and Attribution Registry (SADAR) specification, developed by Cognita AI, is designed to be that infrastructure layer — open, standards-based, and designed for neutral stewardship.*

## **The Path Forward**

Trustworthy AI governance for agentic systems requires action at three levels simultaneously.

At the organizational level, enterprises need to establish the IRB-equivalent governance function before — not after — material agentic AI deployments. This means a defined review process, explicit risk appetite criteria for autonomous AI, disclosure policies for AI-influenced decisions, and the organizational authority to enforce them. It means defining, in advance, what a QA function for AI looks like — what gets sampled, what constitutes an acceptable outcome, what triggers escalation, and to whom. This is governance work, not technology work, and it cannot be delegated to the technology team.

At the standards level, the industry needs an open, vendor-neutral specification for the technical infrastructure that supports governed agentic AI deployment — covering identity, semantic grounding, discovery, process definition, and the cross-boundary audit trail. That specification needs to be developed in collaboration with the standards bodies, regulatory agencies, and industry groups that will ultimately be required to recognize it. The window for establishing these standards before proprietary solutions entrench is open now and will not remain open indefinitely.

At the technical level, organizations deploying agentic AI today can begin building the infrastructure foundation incrementally, without waiting for all five layers of the framework to be fully mature. Starting with verifiable identity and semantic grounding for existing agent components provides immediate governance value and creates the foundation for the higher layers as they develop.

*The question is not whether trustworthy AI governance is achievable. It is whether organizations will build it before or after the events that make its absence undeniable.*

The organizations that move deliberately on this — that establish the governance function, engage with the emerging standards, and build the technical infrastructure as a parallel

workstream to capability deployment — will have a durable advantage. Not only in regulatory posture and defensibility, but in the ability to grant agents the level of autonomy that actually transforms business operations. Governance infrastructure is not a constraint on agentic AI capability. It is the precondition for it.

---

#### **ABOUT THE AUTHOR**

*Clay House is the founder and CEO of Cognita AI Inc. and creator of the CogniWeave™ platform and SADAR specification. He brings a background as CISO at a major health insurer, serves on the Maryland Cybersecurity Council (critical infrastructure subcommittee), and consults on cybersecurity legislation. He has advised on AI governance in clinical and regulated contexts including a panel presentation at HIMSS on securing AI in healthcare.*

---

© 2026 Cognita AI Inc. · CogniWeave™ and SADAR™ are unregistered trademarks of Cognita AI Inc.

*This paper represents the views of the author and is intended to contribute to the public conversation on AI governance.*