

# **Governed Tool Lifecycles for Managed Inference-Provider Ecosystems: A Substrate Approach to Owning, Routing, and Retiring AI Tools**

Teams now assemble applications from dozens of model endpoints and tool providers, but the tools are loose configuration: there is no owner that installs, validates, retrains, and retires them under one auditable policy, and no principled basis for deciding which endpoint should answer a given request. This article shows how a managed tool and inference-provider ecosystem with a governed tool lifecycle and cognitive-state-conditioned dispatch can be built on the Agent-Resident Execution Substrate, disclosed in U.S. Provisional Application No. 64/070,239. It draws on sibling portfolio work in cryptographically enforced agent governance and persistent cognitive state.

---

## **What This Application Specifies**

A managed tool and inference-provider ecosystem is the layer of an AI application where multiple model endpoints and external tools are installed, selected, invoked, updated, and removed over time. In most current deployments this layer is informal: endpoints are named in a configuration file, an orchestration script picks one by hard-

coded rule, and the lifecycle of each tool (who published it, when it was last updated, whether its replacement was validated) lives in tickets and tribal memory rather than in any enforced structure.

This application specifies that ecosystem as a set of managed assets owned by a persistent semantic agent. On the Agent-Resident Execution Substrate, the inference endpoints and external tools are not free-standing services the application happens to call. They are entries in a managed inference tool registry that is subordinate to the agent, and every operation on them runs through a governed tool lifecycle and a cognitive-state-conditioned dispatcher.

Three substrate primitives carry the design. First, the managed inference tool registry stores each endpoint as a record comprising a model artifact, an interface specification, and an associated governance scope, together with a publisher-signed capability declaration that enumerates input modalities, task categories, input and output formats, a latency profile, a model version and publisher identifier, a resource envelope, and governance-compatibility tags. Endpoints of different types and sizes (general-purpose language models, task-specific fine-tuned models, classifiers, embedding and retrieval models, personal corpus models) co-reside in one registry, including adapter-based variants that share a base artifact and apply per-endpoint adapter weights.

Second, the agent-to-tool dispatcher routes each inference request conditioned on the current cognitive state field, the input characteristics of the request, the applicable policy objects, and the per-endpoint capability declaration. The dispatcher derives an input modality and task category, matches them against capability declarations to find candidate endpoints, and consults cognitive domain fields (an execution confidence domain, a normative integrity domain, and a capability awareness domain) to set applicable confidence, normative, and capability thresholds. Candidates falling below those thresholds are excluded, and among the rest a selection function prefers endpoints with higher historical outcome quality recorded in the lineage field for similar requests.

Third, the tool lifecycle controller performs installation, retraining, replacement, archival, removal, quiescence, and resumption as governed operations. Each managed endpoint moves through an explicit state machine (uninstalled, installed-inactive, active, retraining, updated-active, archived, removed) whose transitions are admitted by policy objects and recorded in the agent's append-only lineage field. Retraining and substitution run in a staging area distinct from the active registry, and an updated artifact is promoted only on successful completion and successful policy validation; a failed validation rolls back to the prior artifact and records the cause.

## **Why It Matters**

The market problem is governance debt in multi-provider AI stacks. As an application grows from one model to a portfolio of specialized endpoints and third-party tools, the questions that matter stop being "does the model work" and become operational: which provider answered this request and why, was this endpoint's last update validated before it went live, who published the tool we just installed, and what happens to in-flight traffic when we swap a model. Conventional stacks answer these with external glue: a gateway here, a config registry there, an audit log somewhere else, none of which is authoritative and none of which is bound to a single accountable owner.

The substrate collapses that glue into owned structure. Because the registry, the dispatcher, and the lifecycle controller are all subordinate to one persistent agent, and because every dispatch, lifecycle transition, installation, and validation appends a deterministic record to a single append-only lineage field under a continuity proof, the ecosystem becomes auditable by construction. The complete history of which endpoint was registered when, under which policy version, and with what downstream outcome is deterministically reconstructible from the lineage rather than reassembled from scattered logs.

It also matters for tool sourcing. The spec describes a tool discovery and installation subsystem that retrieves candidate endpoints from declared source authorities (signed publisher repositories, organizational repositories, peer devices, user archives), each candidate arriving with a publisher signature, a capability declaration, and a publisher metadata record. An installation policy object decides admissibility against admissible publisher authorities, capability declarations, artifact size limits, and governance scopes, and rejected candidates are recorded as denial events. Publisher attribution is preserved across lifecycle operations and is available for compensation routing, so the ecosystem can attribute dispatch, retraining, and consultation events back to the provider that supplied the endpoint and compute obligations under a compensation policy object.

## **How It Composes With the Domain**

A managed inference-provider ecosystem built on the substrate composes as follows.

Onboarding a provider is a governed installation. A candidate endpoint is retrieved from a declared source authority with its publisher signature and capability declaration. The installation policy object evaluates it; on admission, the endpoint enters the registry in installed-inactive state through a governed installation lifecycle operation, with the publisher signature and metadata preserved in the record. The capability declaration is verified at installation and at each subsequent lifecycle operation affecting the endpoint. Activation moves it to active state and makes it eligible for dispatch.

Request routing is cognitive-state-conditioned dispatch. When a host application submits an inference request through the application interface (each request carrying a request descriptor, an output requirement, policy hints, an originating host application identifier, and an originating scope context), the dispatcher selects among co-resident endpoints by modality, task category, policy, and capability declaration, gated by the cognitive thresholds described above. The same mechanism supports composition: a

single request can fan out to several endpoints concurrently with outputs aggregated by majority voting, weighted averaging by historical outcome quality, or confidence-weighted selection; or it can run serially through a configured pipeline (for example, a task-classification endpoint, then a task-specific endpoint, then a verification endpoint), reconfigurable through registry updates without changing the endpoints themselves.

Provider updates are atomic, validated substitutions. When a provider ships a new model version, or when the agent's own lineage analysis detects an outcome-quality threshold breach, a begin-retrain or replacement transition stages the new artifact, validates it against policy, and promotes it only on success, with traffic suspended or routed to the prior version during the transition under policy. A bad update never reaches live dispatch, because promotion is conditioned on validation and rolls back on failure. Critically, none of this touches the agent's identity: under the continuity guarantee, the agent's identity, cognitive state, and lineage are preserved across replacement, retraining, or removal of any endpoint, so the ecosystem can replace every endpoint in the registry without disturbing the accountable owner.

Capacity and capability gaps are handled by governed forwarding rather than silent fallback. When no local endpoint satisfies a request's capability requirements or local compute is insufficient, a cloud-burst forwarding subsystem can forward the request to a remote inference endpoint, but only after an admissibility test (a capability test, a capacity test, a disclosure test, and a cost test) and only under a cloud-burst policy object that names admissible remote endpoints, admissible request categories, input disclosure scopes, encryption requirements, and per-request cost limits. Forwarded payloads are treated as off-device disclosure events under the privacy invariant. A confidential-execution forwarding mode encrypts requests under keys held only inside the remote endpoint's trusted execution environment, so payloads are not exposed to the remote operator.

Resource arbitration keeps the ecosystem within budget. A resource governance subsystem enforces per-tool budgets (memory, storage, inference compute, retraining compute, network bandwidth, power), schedules concurrent operations so foreground dispatch is not delayed by background retraining or ingestion, evicts or quiesces endpoints under resource pressure, and records each governance decision in the lineage field.

## **What This Enables**

For platform teams, the substrate enables a provider-agnostic ecosystem with a single accountable owner. Endpoints from many publishers can be installed, ranked, composed, and retired under one governance policy, with routing decisions justified by capability declarations and historical outcome quality rather than by hard-coded rules. Embodiments include a registry of co-resident specialized endpoints (a code-generation endpoint, a writing-assistance endpoint, a domain-specific question-answering endpoint, a media-classification endpoint) each with its own sub-lineage, retraining schedule, and governance scope.

For provider marketplaces, it enables signed, attributable distribution. Because each endpoint carries a verified publisher signature and preserved attribution metadata, the substrate supports a marketplace in which providers publish endpoints to declared source authorities, installations are admitted by policy, and dispatch, retraining, and consultation events are attributed to the publisher for compensation routing under a compensation policy object that declares per-publisher rate schedules and payment destinations.

For regulated and privacy-sensitive deployments, it enables verifiable disclosure control. Lineage records, model artifacts, training corpora, and personal corpus model parameters are not transmitted off the device except under an explicit disclosure policy object naming a recipient, a permitted scope, an authorization attestation, a retention requirement, and a revocation mechanism, and every off-device disclosure (including

federation and cloud-burst events) is recorded as a deterministic, auditable disclosure event. Enforcement options include a substrate-runtime egress filter, per-component isolation, signed disclosure preconditions gating release of transmission keys, and hardware-anchored attestation.

For multi-context operators, scope partitioning enables one agent identity to run independent provider ecosystems per context (professional, personal, project-specific, regulatory-domain-specific), each with its own scope-local tool subset, corpus policy, and lineage partition, while the agent's overall continuity proof spans the unified lineage. Federation extends this across a device population: endpoints are managed locally on each device and lineage records are exchanged under a federation policy, so outcome signals observed on one device improve routing and corpus assembly on another without exchanging model weights.

## **Boundary Conditions**

This is an early-stage disclosure. U.S. Provisional Application No. 64/070,239 establishes priority for the substrate architecture; it is not an issued patent, and the scope of any resulting claims will be determined through examination. The article describes architecture and mechanism, not a benchmarked product, and it asserts no performance figures, latency numbers, throughput, or cost results.

Several boundary conditions follow from the spec itself. The substrate is designed for devices with bounded local memory, storage, and compute, so the registry holds artifacts of bounded size and uses parameter-efficient fine-tuning and quantization; very large frontier models that exceed the local envelope are reached through governed cloud-burst forwarding, not resident execution. Cognitive-state-conditioned dispatch presupposes a populated cognitive state field with calibrated domain thresholds and capability declarations; routing quality depends on the fidelity of those declarations and on accumulated lineage, so a cold-start ecosystem with little outcome history routes on declared capability before outcome-quality preference becomes informative. Publisher

signing and compensation routing assume trustworthy source authorities and signature verification; the substrate governs admission and attribution but does not by itself certify a publisher's claims. Federation and cloud-burst both cross the device boundary and are therefore constrained by the disclosure policy, which may legitimately forbid forwarding for sensitive request categories.

The underlying components are largely prior art and are treated as such: language model inference, low-rank and adapter-based fine-tuning, hardware security elements, trusted execution environments, and content-addressed transport are established techniques. The contribution is the governed composition: a persistent agent that owns its inference endpoints as managed assets, conditions dispatch on its own cognitive state, and binds the entire tool lifecycle to one append-only, continuity-proofed record.

## **Disclosure Scope**

The technology described here (the managed inference tool registry, the cognitive-state-conditioned dispatcher, the governed tool lifecycle and its state machine, atomic validated substitution, publisher signing and compensation routing, cloud-burst forwarding, the resource governance subsystem, and the privacy invariant on off-device disclosure) is disclosed in U.S. Provisional Application No. 64/070,239, "Agent-Resident Execution Substrate." Sibling portfolio applications incorporated by reference in that filing supply the persistent cognitive state representation, the cryptographically signed governance policy objects, and the hardware-anchored identity primitives referenced above; this article does not restate their separate disclosures. The domain framing in this article (managed tool and inference-provider ecosystems, provider marketplaces, multi-provider AI stacks, and any regulatory or compliance context) is external context offered to illustrate a faithful enabling implementation. It is not part of the patent claims, does not characterize the legal scope of any claim, and should not be read as a representation about specific products, providers, standards, or regulatory

obligations. Any reference to publisher compensation, marketplaces, or disclosure auditing describes substrate mechanisms as disclosed, not commercial or legal commitments.

---

## **Agent-Resident Execution**

[All 40 steps → \(/inventive-steps\)](#)

### **Substrate** ([/agent-resident-execution-substrate](#))

Persistent execution environment carried by the agent, not the host — identity, state, and lineage across power cycles, devices, and upgrades.

Provisional application

### **PRIMARY TECHNICAL DISCLOSURE**

- [Agent-Resident Execution Substrate, Articles \(/articles/agent-resident-execution-substrate\)](#)

### **SECONDARY TECHNICAL**

- [Persistent Semantic Agent \(/articles/agent-resident-execution-substrate/persistent-semantic-agent\)](#)
- [Managed Inference Tool Registry \(/articles/agent-resident-execution-substrate/managed-inference-tool-registry\)](#)
- [Agent-to-Tool Dispatcher \(/articles/agent-resident-execution-substrate/agent-to-tool-dispatcher\)](#)
- [Lineage-Derived Training Signal \(/articles/agent-resident-execution-substrate/lineage-derived-training-signal\)](#)
- [Identity Preservation Across Upgrades \(/articles/agent-resident-execution-substrate/identity-preservation-across-upgrades\)](#)
- [Cognitive State-Conditioned Dispatch \(/articles/agent-resident-execution-substrate/cognitive-state-conditioned-dispatch\)](#)
- [Governed Tool Lifecycle \(/articles/agent-resident-execution-substrate/governed-tool-lifecycle\)](#)
- [Continuity-Proof Lineage \(/articles/agent-resident-execution-substrate/continuity-proof-lineage\)](#)
- [Substrate Runtime Continuity \(/articles/agent-resident-execution-substrate/substrate-runtime-continuity\)](#)
- [Personal Corpus Model Training \(/articles/agent-resident-execution-substrate/personal-corpus-model-training\)](#)

- [Heterogeneous Inference Endpoints \(/articles/agent-resident-execution-substrate/heterogeneous-inference-endpoints\)](/articles/agent-resident-execution-substrate/heterogeneous-inference-endpoints).
- [Atomic Lifecycle Substitution \(/articles/agent-resident-execution-substrate/atomic-lifecycle-substitution\)](/articles/agent-resident-execution-substrate/atomic-lifecycle-substitution).
- [Integrity Signal Feedback \(/articles/agent-resident-execution-substrate/integrity-signal-feedback\)](/articles/agent-resident-execution-substrate/integrity-signal-feedback)
- [Hardware-Bound Identity \(/articles/agent-resident-execution-substrate/hardware-bound-identity\)](/articles/agent-resident-execution-substrate/hardware-bound-identity)
- [Cognitive State Append-Only Invariant \(/articles/agent-resident-execution-substrate/cognitive-state-append-only-invariant\)](/articles/agent-resident-execution-substrate/cognitive-state-append-only-invariant).
- [Counterparty Identity Records \(/articles/agent-resident-execution-substrate/counterparty-identity-records\)](/articles/agent-resident-execution-substrate/counterparty-identity-records)
- [Privacy Egress-Controlled Disclosure \(/articles/agent-resident-execution-substrate/privacy-egress-controlled-disclosure\)](/articles/agent-resident-execution-substrate/privacy-egress-controlled-disclosure).
- [Federated Cross-Device Agent Identity \(/articles/agent-resident-execution-substrate/federated-cross-device-agent-identity\)](/articles/agent-resident-execution-substrate/federated-cross-device-agent-identity).

## APPLICATIONS · GENERAL

- [Personal AI Agents That Survive Device Loss: One Continuous Identity and a Private Corpus Across Every Device \(/articles/agent-resident-execution-substrate/personal-cross-device-agents\)](/articles/agent-resident-execution-substrate/personal-cross-device-agents).
- [Enterprise Agent Fleets: Stable Agent Identity and Governed Tool Access Across Model Upgrades and Infrastructure Migration \(/articles/agent-resident-execution-substrate/enterprise-agent-fleets\)](/articles/agent-resident-execution-substrate/enterprise-agent-fleets).
- [Audit-Grade Agent Identity for Regulated Finance and Healthcare: Continuity-Proof Lineage Across the Agent Lifecycle \(/articles/agent-resident-execution-substrate/regulated-industry-agents\)](/articles/agent-resident-execution-substrate/regulated-industry-agents)
- [Edge and On-Device Agents: Hardware-Bound Identity Across Heterogeneous Inference Endpoints \(/articles/agent-resident-execution-substrate/edge-and-on-device-agents\)](/articles/agent-resident-execution-substrate/edge-and-on-device-agents).
- [Agent-to-Agent Commerce With Counterparty Identity Records and Egress-Controlled Disclosure \(/articles/agent-resident-execution-substrate/agent-to-agent-commerce\)](/articles/agent-resident-execution-substrate/agent-to-agent-commerce).
- **[Governed Tool Lifecycles for Managed Inference-Provider Ecosystems: A Substrate Approach to Owning, Routing, and Retiring AI Tools \(/articles/agent-resident-execution-substrate/managed-tool-ecosystems\)](/articles/agent-resident-execution-substrate/managed-tool-ecosystems)**.
- [Proving Unbroken Continuity in Long-Lived Autonomous Systems Across Substrate Migration and Atomic Model Substitution \(/articles/agent-resident-execution-substrate/long-lived-autonomous-systems\)](/articles/agent-resident-execution-substrate/long-lived-autonomous-systems).

---

[Agent-Resident Execution Substrate overview → \(/agent-resident-execution-substrate\)](/agent-resident-execution-substrate)

