

LLM-as-Bootstrap: Why Anchor Inference Engines Shrink as the Lineage Matures

A semantic substrate does not just absorb a large language model, it grows out of one. As the lineage at a given anchor accumulates and its decisions stabilize, the resident inference engine becomes structurally substitutable for a simpler, cheaper mechanism that reproduces those decisions over the observed distribution.

Mechanism

Conventional agentic systems treat the large language model as a permanent runtime cost: every step of reasoning re-invokes a general-purpose model, at full price, no matter how settled the local decision has become. The semantic discovery substrate treats the model differently. Each anchor in the adaptive index hosts a local inference engine that scores and selects among candidate transitions, and the substrate is model-agnostic by construction: the inference engine is a pluggable proposal generator whose outputs are subject to independent admissibility evaluation by the execution substrate. Maturation-driven engine substitution adds a second axis of evolution on top of that property. As traversal lineage accumulates at a stable, well-traversed anchor, the resident inference engine is progressively replaced by a smaller, cheaper mechanism that reproduces the engine's decisions over the observed distribution of transition queries.

The substitution is possible because the parent architecture already records, for every completed traversal, which transitions were admitted, which were rejected, and which were decomposed, along with the semantic mutations and confidence updates at each step. The aggregate of these lineage records at a single anchor is a behavioral specification of that anchor's inference engine over the queries it has actually seen. When that specification becomes structurally evident, continued reliance on a general-purpose model to re-derive it is waste. The substrate compiles a substitute from the lineage and, once the substitute is shown to agree with the resident engine, promotes it to the active engine of the anchor. The substrate therefore becomes progressively cheaper to operate as it matures, without retraining and without any change to the traversal protocol, the discovery object schema, or the governance infrastructure.

Substitution is itself a governed event, not a silent optimization. The decision to substitute, the candidate compiled, the evaluation evidence, the commitment, and any later reversion are each recorded in the anchor's operational lineage, so the engine that produced any given traversal result is always attributable. A substitute never inherits authority; it inherits only the role of proposal generator. The execution substrate continues to evaluate every proposed transition against the same policy, lineage-continuity, entropy, and temporal-validity constraints, so the governance integrity of a traversal is identical whether the proposal came from a large language model or from a lookup table compiled out of that model's own recorded decisions.

Operating Parameters

Each anchor maintains a maturation metric over its lineage. The metric comprises at least three measured quantities: a query-distribution coverage measure, capturing how completely the accumulated lineage spans the space of queries the anchor receives; a decision-stability measure, capturing how consistently the resident engine resolves equivalent queries to equivalent transitions over time; and a substitutability score, estimating how closely a simpler mechanism could reproduce the engine's decisions on

the observed distribution. The metric is local to the anchor and is recomputed as new lineage accrues, so maturation is a property each anchor reaches on its own schedule rather than a global state imposed on the index.

When the maturation metric crosses a policy-defined substitution threshold, the anchor compiles a candidate substitute from its lineage and shadow-evaluates it against the resident engine over a rolling window: both mechanisms see live transitions, the resident engine's decision governs, and the substitute's agreement is measured without affecting traversal outcomes. When agreement over the window exceeds a policy-defined substitution-commitment threshold, the substitute is committed as the active engine. Inbound transitions are then attempted against the substitute first. A transition whose query falls outside the substitute's coverage envelope, or whose substitute-assigned confidence falls below a fallback threshold, is escalated to the original engine or to a higher-capacity engine available on the network. Telemetry indicating substitution-induced degradation, such as a rise in downstream rejections or a drift in resolution rate, triggers reversion to the prior engine, recorded as a governed event in the same lineage.

These parameters expose the cost-quality tradeoff as an operating policy rather than an architectural commitment. A deployment that prioritizes precision sets a high commitment threshold and a wide fallback envelope, substituting conservatively and escalating readily. A deployment whose binding constraint is per-inference cost sets a lower commitment threshold and tolerates a narrower envelope, accepting more aggressive substitution in exchange for a lower steady-state operating cost. The thresholds are themselves governed values, carried in the anchor's policy configuration, so the substitution behavior of an anchor is auditable and bounded rather than emergent.

Substitute Classes

The mechanism compiled in place of a matured engine is selected by the content profile of the anchor's semantic neighborhood. Four classes of substitute are disclosed. A frozen-weight lookup table compiled directly from lineage statistics suits anchors over stable factual content, where the admissible transition for a given query has become effectively fixed: taxonomic structure, citation graphs, biographical and reference data. A small distilled model trained on the anchor's accumulated proposal-outcome pairs suits anchors over moderate-entropy domains such as technical documentation and structured corpora, where the decision boundary is learnable but no longer requires a frontier model to represent. A rule-based matcher derived from frequent admissibility patterns suits anchors over highly structured content governed by an explicit ontology or controlled vocabulary, such as legal citation and regulated terminology. An embedding-similarity scorer tuned to the anchor's neighborhood suits anchors whose admissibility decisions have reduced to a stable similarity threshold.

In every case the substitute reproduces a decision function the resident engine has already demonstrated over the anchor's observed distribution; it does not attempt to generalize beyond that distribution. Queries outside the demonstrated distribution are precisely the queries the coverage envelope routes back to a fuller engine. The substrate thus narrows the role of the expensive model to the frontier of each anchor's experience, while the settled interior of that experience is served by a mechanism whose marginal cost approaches that of a table lookup.

Alternative Embodiments

The mechanism admits embodiments differentiated by deployment economics. In a sovereign or air-gapped embodiment, where a frontier model may be unavailable or prohibitively expensive to host inside the enclave, anchors mature toward compiled substitutes that execute entirely on local hardware, and the coverage envelope routes only genuinely novel queries to a constrained, locally hosted fallback. In an edge

embodiment, anchors resident on resource-limited devices substitute toward lookup tables and rule sets whose memory and compute footprint fit the device, with escalation to a regional engine reserved for queries outside the local envelope.

A tiered-fallback embodiment maintains an ordered chain of engines behind a matured anchor, escalating from the compiled substitute to a small distilled model to a frontier model as successive coverage envelopes are exceeded, so that operating cost scales with the genuine novelty of each query rather than with traffic volume. A migration embodiment runs an incumbent engine and a freshly compiled substitute in parallel during a transition window, promoting the substitute only when measured agreement on live traffic meets the published threshold, and retaining the incumbent as the standing fallback. A federated embodiment lets anchors maintained by independent organizations mature on their own lineage and substitute under their own policy, while the discovery object traverses across them unchanged, because the engine class at each anchor is local and the discovery object's carried semantic state, maintained by the execution substrate rather than the engine, provides continuity across anchor boundaries.

Composition

Maturation-driven substitution composes with the rest of the discovery substrate along three load-bearing seams. First, it composes with the structural separation between proposal and authority. Because the inference engine proposes and the execution substrate disposes, the engine need not be trusted, and replacing it changes nothing about the governance guarantee: admissibility is evaluated identically before and after substitution. This is what makes runtime engine replacement safe at all, and it is a direct consequence of the proposal-authority separation disclosed for the three-in-one traversal step.

Second, it composes with lineage-driven index evolution. The same completed-traversal lineage records that drive anchor splitting, merging, and re-publication also drive engine substitution; the substrate already accumulates the signal, and substitution is a second consumer of it. An anchor that is reinforced by high-value traversals and is therefore preserved by the self-organization mechanism is precisely an anchor whose stable lineage makes its engine a strong substitution candidate. Third, it composes with provenance and confidence: every traversal result already carries an admissibility audit trail, and the substitution layer extends that trail with the identity, version, and evaluation evidence of the engine that produced each proposal, so a downstream consumer can weight or re-examine results by the maturity of the engine behind them.

Prior-Art Distinction

Knowledge distillation trains a smaller model to imitate a larger one, but it operates at training time, against a curated corpus, as a deliberate engineering project, not as a runtime mechanism driven by an anchor's own accumulated decisions and gated by deterministic admissibility. Mixture-of-experts routing selects among several co-resident models for a given input, but it routes among engines rather than replacing one engine with a structurally simpler mechanism compiled from that engine's recorded behavior, and the population of experts is fixed in advance. Small language models, retrieval-augmented small models, and model-cascade approaches reduce cost by choosing a cheaper model up front or by escalating on uncertainty, but none compiles the cheaper mechanism from the system's own operational history, and none ties the substitution to a governed maturation metric recorded in an auditable lineage. The distinguishing combination disclosed here is substitution driven by anchor-local lineage, gated by a shadow-evaluation and commitment protocol, bounded by a coverage envelope with deterministic escalation, reversible on telemetry, and recorded as a governed event, executing beneath an admissibility layer whose guarantee is independent of the engine in place.

Disclosure Scope

The model-agnostic inference primitive, in which any computational mechanism capable of ordering structured candidates may serve as the inference engine at an anchor because the execution substrate provides the governance guarantee independently of the engine, is disclosed in the cognition filing (U.S. Application No. 19/647,395 and its international counterpart) at Section 10.11. The use of completed-traversal lineage as an evolution signal that drives anchor self-organization is disclosed in the same filing at Section 10.20. This article discloses, as an extension of those primitives, the maturation-driven engine-substitution mechanism: the per-anchor maturation metric over lineage; the compilation of a substitute inference mechanism from accumulated proposal-outcome history; the shadow-evaluation and commitment protocol; the coverage envelope with deterministic escalation to a higher-capacity engine; reversion on telemetry indicating degradation; and the recording of substitution and reversion as governed events in the anchor's operational lineage.

The disclosure covers the substitute classes enumerated above (lineage-compiled lookup tables, distilled small models, rule-based matchers, and embedding-similarity scorers), the embodiments differentiated by deployment economics (sovereign, edge, tiered-fallback, migration, and federated), and the compositional seams with proposal-authority separation, lineage-driven index evolution, and provenance. The scope extends to substitute mechanism classes not enumerated whose decision function is compiled from anchor-local lineage and served beneath the same admissibility layer, and to maturation and commitment policies not described whose behavior reduces to the governed-substitution protocol disclosed here.

PRIMARY TECHNICAL DISCLOSURE

- [Governed Semantic Discovery: Search, Inference, and Execution Through Adaptive Traversal \(/articles/governed-semantic-discovery-search-inference-and-execution-through-adaptive-traversal\)](/articles/governed-semantic-discovery-search-inference-and-execution-through-adaptive-traversal)

SECONDARY TECHNICAL

- [The Adaptive Index as Unified Search-Inference-Execution Substrate \(/articles/semantic-discovery/unified-substrate\)](/articles/semantic-discovery/unified-substrate)
- [Three-in-One Traversal: Search, Inference, and Execution in a Single Step \(/articles/semantic-discovery/three-in-one-traversal\)](/articles/semantic-discovery/three-in-one-traversal)
- [The Discovery Object: A Traversal-Native Semantic Agent \(/articles/semantic-discovery/discovery-object\)](/articles/semantic-discovery/discovery-object)
- [Post-PageRank Semantic Ranking: Relevance Through Governed Traversal \(/articles/semantic-discovery/post-pagerank\)](/articles/semantic-discovery/post-pagerank)
- [Persistent Semantic State: Eliminating Prompt Reconstruction \(/articles/semantic-discovery/persistent-state\)](/articles/semantic-discovery/persistent-state)
- [Traversal Lineage as Index Evolution Signal \(/articles/semantic-discovery/traversal-lineage\)](/articles/semantic-discovery/traversal-lineage)
- [Anchor Semantic Neighborhood Publication \(/articles/semantic-discovery/semantic-neighborhoods\)](/articles/semantic-discovery/semantic-neighborhoods)
- [Inference-Time Execution Control as Traversal Primitive \(/articles/semantic-discovery/inference-governance\)](/articles/semantic-discovery/inference-governance)
- [Anchor Self-Organization Under Entropy and Load Pressure \(/articles/semantic-discovery/anchor-self-organization\)](/articles/semantic-discovery/anchor-self-organization)
- [Alias Resolution as Navigational Traversal \(/articles/semantic-discovery/alias-resolution\)](/articles/semantic-discovery/alias-resolution)
- [Three Discovery Operating Modes: Human Search, Agent Reasoning, Answer Synthesis \(/articles/semantic-discovery/operating-modes\)](/articles/semantic-discovery/operating-modes)
- [Model-Agnostic Semantic Discovery \(/articles/semantic-discovery/model-agnostic\)](/articles/semantic-discovery/model-agnostic)
- [Affect-Modulated Discovery Traversal \(/articles/semantic-discovery/affect-modulated-traversal\)](/articles/semantic-discovery/affect-modulated-traversal)
- [Confidence-Gated Discovery Traversal \(/articles/semantic-discovery/confidence-gated-traversal\)](/articles/semantic-discovery/confidence-gated-traversal)
- [Integrity-Tracked Traversal Drift Detection \(/articles/semantic-discovery/integrity-tracked-drift\)](/articles/semantic-discovery/integrity-tracked-drift)
- [Biological Identity-Scoped Access During Discovery \(/articles/semantic-discovery/biological-access\)](/articles/semantic-discovery/biological-access)
- [Rights-Grade Anchor Governance for Content Discovery \(/articles/semantic-discovery/rights-grade-anchors\)](/articles/semantic-discovery/rights-grade-anchors)
- [Forecasting-Shaped Discovery Traversal \(/articles/semantic-discovery/forecasting-shaped\)](/articles/semantic-discovery/forecasting-shaped)
- [Capability-Constrained Anchor Accessibility \(/articles/semantic-discovery/capability-constrained\)](/articles/semantic-discovery/capability-constrained)

- [Collaborative Multi-Object Discovery Traversal \(/articles/semantic-discovery/collaborative-traversal/\)](/articles/semantic-discovery/collaborative-traversal/).
- [Discovery-Driven Sensor Invocation Closed Loop \(/articles/semantic-discovery/sensor-invocation-loop\)](/articles/semantic-discovery/sensor-invocation-loop/).
- [Cross-Platform Credentialed Reader Activation \(/articles/semantic-discovery/credentialed-reader-activation\)](/articles/semantic-discovery/credentialed-reader-activation/).
- **[LLM-as-Bootstrap: Why Anchor Inference Engines Shrink as the Lineage Matures \(/articles/semantic-discovery/maturation-engine-substitution\)](/articles/semantic-discovery/maturation-engine-substitution/)**
- [Personal Cognitive Asset: How Per-User Lineage Re-Weights the Same Substrate \(/articles/semantic-discovery/personal-lineage-layer\)](/articles/semantic-discovery/personal-lineage-layer/).
- [Loki, the Dog, and the Symbol Grounding Problem \(/articles/semantic-discovery/hybrid-symbol-grounding\)](/articles/semantic-discovery/hybrid-symbol-grounding/)

APPLICATIONS · GENERAL

- [Enterprise Knowledge Management Through Governed Traversal \(/articles/semantic-discovery/enterprise-knowledge-management\)](/articles/semantic-discovery/enterprise-knowledge-management/).
- [AI-Native Search That Replaces PageRank With Contextual Relevance \(/articles/semantic-discovery/ai-native-search\)](/articles/semantic-discovery/ai-native-search/).
- [Semantic Discovery for Scientific Research \(/articles/semantic-discovery/scientific-research-discovery\)](/articles/semantic-discovery/scientific-research-discovery/).
- [Semantic Discovery for Legal Case Research \(/articles/semantic-discovery/legal-case-research\)](/articles/semantic-discovery/legal-case-research/).
- [Semantic Discovery for Patent Landscape Analysis \(/articles/semantic-discovery/patent-landscape-analysis\)](/articles/semantic-discovery/patent-landscape-analysis/).
- [Semantic Discovery for Medical Literature Search \(/articles/semantic-discovery/medical-literature-search\)](/articles/semantic-discovery/medical-literature-search/).
- [Semantic Discovery for Competitive Intelligence \(/articles/semantic-discovery/competitive-intelligence\)](/articles/semantic-discovery/competitive-intelligence/).
- [Semantic Discovery for Regulatory Compliance Search \(/articles/semantic-discovery/regulatory-compliance-search\)](/articles/semantic-discovery/regulatory-compliance-search/).
- [Discovery-Coordinated Multi-Sensor Perception \(/articles/semantic-discovery/coordinated-perception\)](/articles/semantic-discovery/coordinated-perception/).
- [Post-AirTag Cross-Platform Object Tracking \(/articles/semantic-discovery/post-airtag-tracking\)](/articles/semantic-discovery/post-airtag-tracking/).
- [Use the World as Memory: The Brain Strategy for AI \(/articles/semantic-discovery/world-as-memory\)](/articles/semantic-discovery/world-as-memory/)

APPLICATIONS · SPECIFIC

- [Google Search Retrieves Results, Not Understanding \(/articles/semantic-discovery/google-search\)](/articles/semantic-discovery/google-search)
- [Perplexity Answers Questions Without Discovery State \(/articles/semantic-discovery/perplexity\)](/articles/semantic-discovery/perplexity)
- [Elasticsearch Indexes Documents, Not Discovery \(/articles/semantic-discovery/elasticsearch\)](/articles/semantic-discovery/elasticsearch)
- [Algolia Optimizes Relevance Without Discovery State \(/articles/semantic-discovery/algolia\)](/articles/semantic-discovery/algolia)
- [Pinecone Finds Vectors, Not Understanding \(/articles/semantic-discovery/pinecone\)](/articles/semantic-discovery/pinecone)
- [Weaviate Stores Semantics Without Discovery Governance \(/articles/semantic-discovery/weaviate\)](/articles/semantic-discovery/weaviate)
- [You.com Answers Questions but Does Not Govern Discovery \(/articles/semantic-discovery/you-com\)](/articles/semantic-discovery/you-com)
- [Brave Search Built an Independent Index Without Governed Traversal \(/articles/semantic-discovery/brave-search\)](/articles/semantic-discovery/brave-search)
- [Kagi Charges for Better Results, Not Governed Discovery \(/articles/semantic-discovery/kagi\)](/articles/semantic-discovery/kagi)
- [Metaphor Systems Predicts Links but Does Not Govern Traversal \(/articles/semantic-discovery/metaphor-systems\)](/articles/semantic-discovery/metaphor-systems)
- [Glean Indexes Enterprise Knowledge Without Governing Its Discovery \(/articles/semantic-discovery/glean\)](/articles/semantic-discovery/glean)
- [Coveo Personalizes Retrieval, Not Discovery Governance \(/articles/semantic-discovery/coveo\)](/articles/semantic-discovery/coveo)
- [Apple Find My Lacks Cross-Authority Reader Activation \(/articles/semantic-discovery/apple-find-my\)](/articles/semantic-discovery/apple-find-my)
- [Google Find My Network Needs Credentialed Cross-Activation \(/articles/semantic-discovery/google-find-my\)](/articles/semantic-discovery/google-find-my)
- [IETF DULT Specifies Behavior, Not Architecture \(/articles/semantic-discovery/ietf-dult\)](/articles/semantic-discovery/ietf-dult)
- [Glean Enterprise Search and Work AI \(/articles/semantic-discovery/glean-enterprise-search\)](/articles/semantic-discovery/glean-enterprise-search)
- [GraphRAG, but with Governance: Where Microsoft's Architecture Stops Short \(/articles/semantic-discovery/microsoft-graphrag\)](/articles/semantic-discovery/microsoft-graphrag)
- [Memory Layers for Agents: Why Mem0, Zep, and Letta Get Close \(/articles/semantic-discovery/memory-for-agents\)](/articles/semantic-discovery/memory-for-agents)

[Semantic Discovery overview → \(/semantic-discovery\)](/semantic-discovery)