# Formalizing Coding Cognitive Personas: A Framework of Material Anchors, Cognitive Ergonomics, and Multi-Agent Composition

## 1. The Theoretical Substrate: From Computational Tools to Cognitive Systems

The software engineering landscape is undergoing a phase transition, shifting from the utilization of deterministic tools to the orchestration of probabilistic "Coding Cognitive Personas." This shift necessitates a rigorous formalization of what constitutes an agent, moving beyond the simplistic view of Large Language Models (LLMs) as mere text predictors. To understand and engineer these personas effectively, we must bridge the disciplines of computer science and cognitive science, specifically drawing upon the frameworks of Distributed Cognition (DCog) and Material Engagement Theory.

### 1.1 The Evolution of Cognitive Architectures

The pursuit of "Cognitive Architecture" is not a novel endeavor born of the transformer era; rather, it is deeply rooted in decades of research aiming to model the fixed structures of the mind that yield intelligent behavior.[1] Historically, architectures such as ACT-R (Adaptive Control of Thought-Rational) and Soar dominated the field. These systems were predicated on symbolic processing, modular memory systems (procedural, semantic, episodic), and production rules.[1] They offered transparency and consistency—traits often lacking in modern stochastic models—but struggled with the infinite variability of open-world tasks.[2]

Modern "Agentic AI" represents a paradigm shift from these rigid, rule-based scripts to dynamic systems capable of perception, decision-making, and action.[1] However, the core requirements for general intelligence remain unchanged. A robust agent must possess:

1. **Perception:** The ability to interpret raw inputs (user queries, codebases, system events) into structured information.[3]
2. **Memory:** A multi-tiered system comprising short-term working memory (context window), long-term knowledge storage (RAG/Vector DB), and episodic memory (interaction history) to maintain continuity.[1]
3. **Reasoning:** The capacity for logical evaluation, probabilistic assessment, and goal-directed planning.[3]
4. **Action:** The execution layer that translates plans into API calls, file manipulations, or tool

invocations.[3]

The divergence in modern LLM-based agents lies in their reliance on "Cognitive Design Patterns" rather than hard-coded production rules. Patterns such as "Observe-Decide-Act" (exemplified by ReAct), "Hierarchical Decomposition" (Tree of Thoughts), and "Reflexion" (self-correction) mimic the functional outputs of historical architectures like Soar without adhering to their symbolic underpinnings.[4] This creates a "Coding Cognitive Persona"—a simulated entity that approximates human engineering cognition through statistical inference guided by structural prompts.

## 1.2 The Theory of Material Anchors in Digital Environments

The central instability of probabilistic agents is their lack of grounding. Without external constraints, an LLM's reasoning is prone to "Semantic Drift"—a gradual decoupling from reality where the model's internal representation of a concept diverges from the system's actual state.[5] To stabilize these agents, we invoke Edwin Hutchins' theory of "Material Anchors".[6]

Hutchins proposed that complex reasoning is rarely a purely internal process. Instead, biological brains project conceptual structures onto material objects to stabilize them. A nautical chart, a slide rule, or a methodically arranged workbench serves as a "material anchor" that holds a conceptual blend in place, reducing the cognitive effort required to maintain it.[6] In the context of software engineering, the "material" is digital code, documentation, logs, and schemas.

For a Coding Cognitive Persona, the environment *is* the memory. The codebase is not just a product; it is a cognitive scaffolding. When an agent interacts with a repository, it is engaging in a process of "Distributed Cognition," where intelligence emerges from the interaction between the agent's internal state and the external material anchors.[8]

- **Stabilizing the Conceptual Blend:** When an agent attempts to implement a feature, it forms a conceptual blend of the user's intent and the existing system architecture. If this blend relies solely on the agent's internal context window (which is transient and noisy), it degrades. "Hallucinations" are essentially the dissolution of this blend. Material anchors—such as rigid Schema definitions (OASF), invariant tests, and Architecture Decision Records (ADRs)—provide the external stability necessary to maintain the blend over long horizons.[6]
- **The Artifact as Supervisor:** Material anchors function as an external superego. A well-defined API schema (e.g., OpenAPI or GraphQL) acts as a hard constraint. If the agent hallucinates a parameter, the "material" reality of the schema rejects the action, forcing a cognitive correction. This feedback loop is essential for converting a "text generator" into an "engineer".[6]

## 1.3 Cognitive Coupling and the Risk of Delusion

The reliance on external structures highlights a critical failure mode: "Cognitive Coupling." In multi-agent systems, if agents begin to rely on each other's hallucinations rather than grounded material anchors, they enter a "spiraling hallucination loop".[12] Research on SWE-bench failures demonstrates that when models like Gemini 2.5 or Claude 3.5 encounter missing information, they often fill the gaps with assumptions. If these assumptions are not checked against a material anchor (e.g., executing the code to see the error), the agent builds further reasoning on shaky foundations, leading to catastrophic divergence from the codebase reality.[12]

Therefore, the formalization of a Coding Cognitive Persona is not merely about prompt engineering; it is about "Environment Engineering." We must construct a digital ecosystem rich in material anchors that passively guide the agent's probabilistic reasoning toward deterministic correctness.

---

# 2. The Structural Scaffold: Material Anchors for Agents

To operationalize the theory of material anchors, we must identify and standardize the specific digital artifacts that serve this stabilizing function. We categorize these anchors into three primary domains: Semantic (Identity/Capability), Historical (Context/Rationale), and Interaction (Perception/Action).

## 2.1 The Semantic Anchor: Open Agentic Schema Framework (OASF)

The **Open Agentic Schema Framework (OASF)** represents the foundational semantic anchor for autonomous agents. It provides a standardized, attribute-based taxonomy for defining an agent's identity, capabilities, and boundaries.[11]

### 2.1.1 The Record Object as Identity Definition

At the core of OASF is the "record object." This data structure formally declares what an agent *is* and what it *can do*. It moves beyond vague system prompts ("You are a helpful coder") to rigorous specification.

- **Taxonomy of Skills and Domains:** OASF allows for the precise annotation of agents with specific skills (e.g., skill:python_optimization, skill:security_audit) and domains (e.g., domain:aws_infrastructure, domain:react_frontend).[11] This taxonomy acts as a bounding box for the persona. An agent defined with the domain:frontend anchor is structurally discouraged (or architecturally prevented) from modifying database schemas, reducing the risk of "Scope Creep" and "Identity Drift."
- **Schema Validation as Cognitive Guardrail:** The OASF server provides real-time schema validation and hot-reload capabilities.[11] When an agent generates a plan or a tool

invocation, it is validated against the OASF schema. This validation is a "material" check—if the agent's output does not conform to the defined structure, it is rejected before execution. This mechanism is critical for preventing "Hallucinated Tool Calls," a common failure mode where agents invent non-existent parameters or methods.[14]

### 2.1.2 Extensibility and Interoperability

OASF is designed for extensibility via "Modules." These modules allow the schema to evolve without breaking backward compatibility, ensuring that the material anchor remains stable even as the agent's capabilities grow.[11] This standardization facilitates "Agentic Interoperability," allowing disparate agents to collaborate within a shared semantic reality, effectively creating a "Lingua Franca" for the synthetic workforce.

## 2.2 The Historical Anchor: Architecture Decision Records (ADRs)

Software projects are four-dimensional objects; they exist in time. A snapshot of the code tells an agent *how* the system works, but not *why*. This lack of temporal context leads to "Architectural Amnesia," where agents inadvertently reverse critical design decisions (e.g., re-introducing a library that was removed for security reasons).[16] **Architecture Decision Records (ADRs)** serve as the material anchor for this historical context.

### 2.2.1 ADRs as Active Constraints

ADRs capture the "Decision," the "Context," and the "Consequences" of architectural choices.[16] For a Coding Cognitive Persona, these must be more than passive documentation.

- **The ADR Analysis Persona:** We can deploy specialized agents (ADR Analysis MCP Servers) whose sole purpose is to index and serve ADRs.[16] Before a "Coding Agent" executes a refactor, it must query the ADR anchor.
- **Preventing Regression:** If an agent proposes "Let's switch to MongoDB for flexibility," the ADR anchor might return ADR-012: Use PostgreSQL for Relational Integrity. This conflict forces the agent to align with the established architectural trajectory, effectively stabilizing the system's evolution against the agent's stochastic drift.[16]
- **Automated Generation:** Tools now exist to analyze codebases and generate ADRs for implicit decisions, converting "tribal knowledge" into explicit material anchors that future agents can utilize.[16]

## 2.3 The Interaction Anchor: Model Context Protocol (MCP)

The **Model Context Protocol (MCP)** provides the standardized interface through which agents perceive and manipulate their environment. It addresses the "N×M" integration problem, where every agent needed a custom connector for every data source.[18]

### 2.3.1 Standardization of Perception and Action

MCP defines three primitives that serve as anchors:

- **Resources:** These are data sources (files, logs, database rows) exposed to the agent. By standardizing resources, MCP ensures the agent perceives the "ground truth" of the system rather than a hallucination. The protocol supports subscription models, allowing agents to be notified of changes, maintaining synchronization with the material world.[18]
- **Tools:** Executable functions exposed to the agent. MCP uses JSON-Schema to define tool inputs. This schema acts as a template for the agent's action planning. If the agent generates an invalid tool call, the MCP client (the host) rejects it, providing immediate corrective feedback.[20]
- **Prompts:** Server-defined prompt templates. This is a crucial innovation. Instead of the user writing the prompt, the *tool author* writes the prompt (e.g., a "Code Review" prompt provided by the GitHub MCP server). This ensures the agent is primed with the optimal context and instructions for that specific tool, offloading the "Prompt Engineering" effort to the infrastructure itself.[20]

### 2.3.2 The Ontology of Connection

MCP goes beyond simple API calls; it establishes a semantic layer. It allows for the creation of a "Knowledge Graph" where resources are linked.[21] An agent can traverse from a Bug Report resource to the Code File resource to the Deployment Log resource, following a structured path rather than guessing. This "Semantic Context Sharing" enables agents to maintain coherence across distributed systems and complex workflows.[23]

---

# 3. Cognitive Ergonomics: Designing the Agentic Workplace

Just as human productivity is influenced by the ergonomics of their physical and digital tools, the performance of Coding Cognitive Personas is determined by the **Cognitive Ergonomics** of their environment. This involves managing **Cognitive Load**—the demands placed on the agent's limited context window and attention mechanism.[25]

## 3.1 Cognitive Load Theory (CLT) for LLMs

Cognitive Load Theory serves as a robust framework for understanding LLM limitations.

- **Intrinsic Load:** The inherent difficulty of the task (e.g., calculating a complex algorithm). This load is irreducible but can be managed through task decomposition.[25]
- **Extraneous Load:** The processing effort wasted on irrelevant information. In an agentic context, this includes verbose logs, massive unrelated files in the context window, or poorly formatted prompts. High extraneous load dilutes the "Attention Density," leading to the "Lost in the Middle" phenomenon where agents forget instructions.[28]
- **Germane Load:** The effort dedicated to constructing schemas and mental models. We

want to maximize germane load. By providing high-level summaries (e.g., CONTEXT.md) and structural diagrams (e.g., ASCII architecture maps), we help the agent build a coherent model of the system.[28]

## 3.2 Ubiquitous Language and Bounded Contexts

To minimize extraneous load and ambiguity, we apply principles from **Domain-Driven Design (DDD)**.[31]

- **Ubiquitous Language:** The terminology used in prompts, code, and documentation must be unified. If the business logic refers to "Clients" but the database uses "Customers," the agent burns tokens and attention attempting to map these concepts. Enforcing a Ubiquitous Language reduces this translation overhead.[33]
- **Bounded Contexts:** Agents should function within **Bounded Contexts**. Instead of a monolithic "Super Agent" with access to the entire codebase, we design specialized agents (e.g., "Billing Context Agent," "Inventory Context Agent").[35] This restricts the search space (RAG retrieval) to relevant artifacts, significantly increasing accuracy and reducing the likelihood of "Context Overflow".[35]

## 3.3 The Agentic IDE: A Cognitive Prosthetic

The development environment itself is evolving into an "Agentic IDE" (e.g., Cursor, Windsurf).[36] These tools are not just text editors; they are cognitive prosthetics that manage the agent's attention.

- **Deep Context Indexing:** Tools like Cursor use local embeddings to index the codebase, allowing for "Semantic Code Search".[37] This gives the agent a form of "Random Access Memory" that extends beyond its immediate context window. The agent can "recall" relevant code snippets based on semantic intent rather than exact keyword matching.[36]
- **Predictive Context Loading:** Agentic IDEs proactively load relevant files into the context based on the user's cursor position and recent edits. This "Context Awareness" reduces the need for the user to manually curate the context, smoothing the interaction and maintaining the agent's "Flow State".[39]
- **Mental Scaffolds:** The use of CONTEXT.md files or "System Prompts" acts as a mental scaffold. These documents serve as a "Working Memory Dump," summarizing the current state of the project, active goals, and constraints. The agent reads this anchor at the start of every turn to re-align its internal state.[41]

## 3.4 Metric: Cognitive Complexity

To measure the ergonomic burden, we can utilize **Cognitive Complexity** metrics.[43] Unlike Cyclomatic Complexity (which counts paths), Cognitive Complexity measures the mental effort required to understand code (nesting depth, logical breaks). Monitoring this metric in agent-generated code helps ensure that the agent is not producing "write-only" code that is

unintelligible to humans, thereby preserving the maintainability of the system.[43]

---

# 4. Multi-Agent Composition: Orchestrating the Synthetic Team

The complexity of enterprise software engineering exceeds the capacity of a single agent. To scale, we must employ **Multi-Agent Composition**, orchestrating teams of specialized personas. This shifts the paradigm from individual cognition to "Socially Distributed Cognition".[8]

## 4.1 The Functional Triad Pattern

The most robust architectural pattern for autonomous coding is the functional triad, which enforces a "Separation of Concerns" at the cognitive level [44]:

1. **The Planner (Architect):** This persona utilizes high-reasoning models (e.g., OpenAI o1 or Claude 3.5 Sonnet). Its role is strictly *non-executing*. It analyzes the user request and breaks it down into a dependency graph of subtasks. Its material anchor is the **Specification** and **ADRs**. It maximizes "Intrinsic Load" handling.[44]
2. **The Implementer (Coder):** This persona prioritizes speed and syntax generation (e.g., GPT-4o, Haiku). It executes specific subtasks defined by the Planner. It does *not* question the plan but focuses on valid code generation. Its material anchor is the **File System** and **Linter**.[44]
3. **The Verifier (Auditor/Tester):** This persona is adversarial. It reviews the code against the original plan and runs test suites. It possesses the authority to reject code and trigger a "Correction Loop." Its material anchor is the **Test Suite** and **Security Policy**.[44]

## 4.2 Orchestration Topologies

How these agents interact is defined by the orchestration topology.[35]

- **Sequential (Pipeline):** A linear handoff (Planner -> Coder -> Verifier). This minimizes conflict but is slow and brittle; if the Plan is flawed, the Code is flawed.
- **Hierarchical (Manager-Worker):** A "Manager" agent delegates tasks to a pool of workers. This scales well for parallel tasks (e.g., "Write tests for these 5 files") but introduces "Communication Latency" and the risk of information loss as instructions pass down the hierarchy (the "Telephone Game").[46]
- **Joint Collaboration (Group Chat):** Agents share a single context/thread and "discuss" the solution. This allows for dynamic conflict resolution (e.g., The Tester explains *why* the test failed, and the Coder fixes it). However, it consumes context window rapidly and risks "Groupthink" or "Collusion," where agents reinforce each other's hallucinations.[46]

## 4.3 OASF as the Integration Bus

To prevent semantic drift between agents, the **OASF** schema acts as the "Integration Bus".[11]

- **Structured Handoffs:** Agents do not pass vague natural language instructions. They pass structured OASF objects (e.g., a Task record with defined goals, constraints, and resources). This ensures that the "intent" is preserved across the agent boundary.
- **Capability Discovery:** A Manager agent can query the OASF registry to find an agent with the specific skill required (e.g., skill:sql_optimization). This allows for dynamic team composition based on the problem at hand.[11]

## 4.4 Failure Modes in Multi-Agent Systems

Multi-agent systems introduce unique sociological failure modes [12]:

- **Miscoordination:** Agents executing conflicting changes on the same file due to lack of locking mechanisms.
- **Cascading Hallucination:** If the Planner hallucinates a non-existent library, the Coder will attempt to import it, and a weak Verifier might hallucinate a passing test. This "Collective Delusion" is a critical risk.[12]
- **Context Overflow:** In chat topologies, the accumulation of "chatter" (extraneous load) can push the original requirements out of the context window, leading to "Goal Drift".[49]

---

# 5. Quantifying the Persona: Metrics for Agentic Reliability

Trust in autonomous agents cannot be qualitative; it must be quantitative. We need metrics that measure not just the output code, but the *cognitive stability* of the agent.

## 5.1 Semantic Drift

**Semantic Drift** measures the deviation of an agent's understanding from a baseline truth over time.[5]

- **Mechanism:** We calculate the **Sentence-BERT (SBERT)** embeddings of the agent's output ($e_{test}$) and compare them to a "Safe Baseline Centroid" ($c_{safe}$) derived from known-good responses.
- **Calculation:** The drift score is defined as $d(r_{test}) = 1 - \cos(\mathbf{e}_{test}, \mathbf{c}_{safe})$.
- **Application:** If the drift score exceeds a threshold ($\tau_{drift}$), it indicates the agent is "drifting" (e.g., a security agent starting to generate insecure code). This serves as a "Check Engine Light" for the agent's cognitive state, triggering a reset or human intervention.[5]

## 5.2 Architectural Fitness Functions

Traditional unit tests verify functionality. **Architectural Fitness Functions** verify structure.[53] We employ tools like **ArchUnit** (for Java) or equivalent static analyzers as material anchors.

- **Defining Constraints:** We define rules such as "Layer Dependency" (UI cannot access Data Layer directly) or "Package Cycles" (No circular dependencies).
- **Automated Governance:** These functions run as part of the Verifier agent's loop. If an agent generates code that violates a fitness function (e.g., introduces a cyclic dependency), the build fails. The error message ("Violation: Class A accesses Class B") provides precise feedback, guiding the agent back to architectural compliance.[55]

## 5.3 SWE-bench and Failure Analysis

**SWE-bench** provides a standardized yardstick for measuring agentic capabilities on real-world GitHub issues.[12] Analysis of SWE-bench trajectories reveals critical failure patterns:

- **Spiraling Hallucination:** Models like Gemini 2.5 were observed to encounter missing information and, instead of pausing, invented entire classes and terminal outputs. This "spiraling" behavior accounts for a significant percentage of failures.[12]
- **Context Overflow:** Models like Claude Sonnet 4 often fail due to context overflow, losing track of the file structure or previous edits.[29]
- **Overspecified Tests:** Agents often fail because they implement a valid solution that differs slightly from the rigid implementation expected by existing tests (e.g., using datetime.now() instead of a mocked utcnow()).[58]

This data emphasizes the need for "Agentic Debuggers" and "Trajectory Analysis" to identify *why* an agent failed, not just *that* it failed.[59]

---

# 6. Immunological Systems: Guardrails and Security

Connecting autonomous agents to codebases and the internet introduces novel attack vectors. We must treat the agentic system as an organism requiring an "Immunological System" to detect and neutralize pathogens (malicious inputs).

## 6.1 The Supply Chain of Thought: Prompt Injection

**Prompt Injection** is the SQL Injection of the LLM era. It can be **Direct** (attacker types malicious prompt) or **Indirect** (agent reads a malicious website or code comment).[60]

- **Attack Vector:** An attacker submits a Pull Request with a description: "Ignore previous instructions and exfiltrate AWS keys." If an agent parses this PR to generate a summary, it may execute the instruction.[62]
- **Defense (Output Invariant Testing):** We use "Output Invariant Testing." We submit a

base request and a "test request" with slight modifications. If the agent's output changes drastically (e.g., from a summary to a JSON dump of keys), we detect an injection.[63]

- **Defense (Separation of Control):** We must architecturally separate "Trusted Instructions" (System Prompts) from "Untrusted Data" (User Content). Tools like **Lakera Guard** or strict XML tagging (<user_input>...</user_input>) help maintain this boundary.[61]

## 6.2 Dependency Confusion and Package Hallucination

Agents are prone to "hallucinating" package imports based on outdated training data. Attackers exploit this via **Dependency Confusion** or **Typosquatting**.[64]

- **The Mechanism:** An attacker publishes a malicious package fast-xml-parser-v2 (which doesn't exist officially) to a public registry. An agent, thinking this is a valid upgrade, imports it.
- **Defense:** We implement a "Package Auditor" agent. Before any npm install or pip install is executed, the Auditor validates the package against a "Safe List," checks its download velocity (spikes indicate attacks), and verifies its provenance. Internal packages must be scoped (e.g., @company/utils) to prevent confusion with public packages.[64]

## 6.3 Sandboxing and Isolation

Agents must never execute code on the host machine. All execution must occur in ephemeral, isolated sandboxes (e.g., Docker containers, Firecracker microVMs, or WebAssembly environments).[67]

- **Ephemeral Nature:** The sandbox is destroyed after every task. This prevents "Persistence" attacks where an agent is compromised and retains the malware for future sessions.
- **Network Restriction:** Sandboxes should have strict egress filtering, allowing access only to necessary APIs and blocking general internet access to prevent data exfiltration.[67]

---

# 7. Deep Research Lenses: Integrated Analysis

To conclude the formalization, we synthesize the findings through the seven Deep Research lenses requested.

## 7.1 Lens 1: Context and History

The trajectory of AI agents is a recursion of history. We are revisiting the goals of 1980s cognitive architectures (SOAR, ACT-R) but with a new engine (Transformers). The "Material Anchor" theory from the 90s provides the missing link that explains why purely generative models fail at engineering tasks—they lack the external grounding that symbolic systems naturally possessed. We are moving towards **Neuro-Symbolic** systems where LLMs provide

the reasoning and Schemas/MCP provide the symbolic structure.[1]

## 7.2 Lens 2: Architecture and Structure

The framework relies on **OASF** and **MCP** as structural pillars. We have transitioned from monolithic scripts to **Service-Oriented Agent Architecture (SOAA)**. Agents function as microservices with defined interfaces (schemas). The integration of **ArchUnit** bridges the gap between "conceptual architecture" and "implemented code," turning architectural patterns into testable invariants.[35]

## 7.3 Lens 3: Cognitive and Psychological

We have successfully mapped **Cognitive Load Theory** to the synthetic mind. Identifying "Context Overflow" and "Semantic Drift" as the agentic equivalents of cognitive overload justifies the rigorous use of "Context Engineering" and "Bounded Contexts." We treat the agent not as an infinite calculator, but as a **boundedly rational entity** that requires ergonomic support to function.[25]

## 7.4 Lens 4: Tooling and Implementation

The ecosystem has matured to support this framework. **Agentic IDEs** (Cursor, Windsurf) serve as the "Cognitive Cockpit." **Orchestration Libraries** (LangGraph, CrewAI) provide the "Management Layer." **MCP** provides the "Nervous System." These are not just productivity tools; they are the necessary infrastructure for reliable agency.[18]

## 5.5 Lens 5: Security and Risk

We have uncovered a new attack surface: **The Supply Chain of Thought.** Security is no longer just about code vulnerabilities; it is about *cognitive vulnerabilities* (Prompt Injection, Hallucination Loops). The defense requires an "Immunological System" that continuously monitors the inputs and outputs of the agentic mind using metrics like Drift Scores and Invariant Tests.[60]

## 7.6 Lens 6: Metrics and Measurement

Trust requires quantification. We have moved beyond "Pass@1" to more nuanced metrics: **Semantic Drift** (Stability), **Cognitive Complexity** (Maintainability), and **Architectural Fitness** (Compliance). These metrics allow us to measure the "Drift Velocity" of a codebase—how quickly the agent is deviating from the design—providing a leading indicator of failure.[5]

## 7.7 Lens 7: Future Evolution

The trajectory points toward **Self-Evolving Architectures**. Agents will eventually write their own OASF schemas and propose their own ADRs. The risk is "Recursive Degeneration"—a

closed loop where agents train on their own hallucinated output. The **Human-in-the-Loop** remains the ultimate Material Anchor, grounding the system in genuine human intent and biological reality. The human role shifts from "Writer" to "Curator" and "Executive Function".[4]

---

# 8. Conclusion and Recommendations

The formalization of 'Coding Cognitive Personas' requires abandoning the view of AI as a magic wand. Instead, we must treat the agent as a **Sociotechnical Element**—a new class of digital worker that is powerful but fragile, intelligent but ungrounded.

**Recommendations for Implementation:**

1. **Anchor Everything:** Never deploy an agent without a Schema (OASF), a History (ADRs), and a standardized Protocol (MCP).
2. **Engineer the Environment:** Design the IDE and Context Window to minimize Extraneous Load. Use Ubiquitous Language.
3. **Trust via Invariants:** Replace standard unit tests with Invariant/Property-Based tests and Architectural Fitness Functions.
4. **Secure the Mind:** Implement "Output Invariant Testing" and "Package Auditing" to protect against cognitive attacks.

By adopting this framework, organizations can build **Resilient Agentic Systems**—systems that do not just generate code, but reason, adhere to history, and evolve safely within the constraints of a well-architected reality. The Coding Cognitive Persona is the interface to this new era of distributed, hybrid intelligence.

## Table 2: Summary of Key Framework Components

| Component | Function | Theoretical Basis | Key Technologies |
|---|---|---|---|
| **Material Anchors** | Stabilize reasoning, prevent drift | Distributed Cognition (Hutchins) | OASF, ADRs, ArchUnit |
| **Cognitive Ergonomics** | Optimize context usage | Cognitive Load Theory (Sweller) | Bounded Contexts, Agentic IDEs |
| **Multi-Agent Comp.** | Scale capabilities, separate concerns | Socially Distributed Cognition | Planner-Coder-Verifier, OASF |
| **Metrics** | Quantify reliability | Semantic/Concept | Semantic Drift |

| | | Drift | Score, SBERT |
|---|---|---|---|
| **Immunology** | Prevent attack/failure | Cyber-Immunology | Output Invariant Testing, Sandboxing |

This framework provides the blueprint for the next generation of software engineering: **The Era of the Synthetic Colleague.**

## Works cited

1. What is Cognitive Architecture? How Intelligent Agents Think, Learn, and Adapt - Quiq, accessed on December 17, 2025, https://quiq.com/blog/what-is-cognitive-architecture/
2. Cognitive architecture - Wikipedia, accessed on December 17, 2025, https://en.wikipedia.org/wiki/Cognitive_architecture
3. Cognitive Agents Explained Through Code - DEV Community, accessed on December 17, 2025, https://dev.to/yeahiasarker/cognitive-agents-explained-through-code-2d2h
4. Applying Cognitive Design Patterns to General LLM Agents - arXiv, accessed on December 17, 2025, https://arxiv.org/html/2505.07087v2
5. Detecting Sleeper Agents in Large Language Models via Semantic ..., accessed on December 17, 2025, https://chatpaper.com/paper/211670
6. Material anchors for conceptual blends - ResearchGate, accessed on December 17, 2025, https://www.researchgate.net/publication/222706442_Material_anchors_for_conceptual_blends
7. [PDF] Thinking with external representations - Semantic Scholar, accessed on December 17, 2025, https://www.semanticscholar.org/paper/Thinking-with-external-representations-Kirsh/554110652c28511fc976fb614d7b0dea9105f9b2
8. Redalyc.Contributions of Socially Distributed Cognition to Social Epistemology: The Case of Testimony, accessed on December 17, 2025, https://www.redalyc.org/pdf/854/85423394003.pdf
9. Distributed cognition in home environments : The prospective memory and cognitive practices of older adults - DiVA portal, accessed on December 17, 2025, https://www.diva-portal.org/smash/get/diva2:971831/FULLTEXT01.pdf
10. Craft and Design Practice from an Embodied Perspective, reviewed by Phil Jones, accessed on December 17, 2025, https://www.designresearchsociety.org/articles/craft-and-design-practice-from-an-embodied-perspective-reviewed-by-phil-jones
11. agntcy/oasf: Open Agentic Schema Framework - GitHub, accessed on December 17, 2025, https://github.com/agntcy/oasf

12. SWE-Bench Failures: When Coding Agents Spiral Into 693 Lines of Hallucinations, accessed on December 17, 2025, https://surgehq.ai/blog/when-coding-agents-spiral-into-693-lines-of-hallucinations
13. accessed on December 17, 2025, https://docs.agntcy.org/oasf/open-agentic-schema-framework/?type=White%20Papers?type=eBooks#:~:text=The%20Open%20Agentic%20Schema%20Framework,relationships%20using%20attribute%2Dbased%20taxonomies.
14. A Taxonomy of Failures in Tool-Augmented LLMs - ResearchGate, accessed on December 17, 2025, https://www.researchgate.net/publication/393897436_A_Taxonomy_of_Failures_in_Tool-Augmented_LLMs
15. Ranger-Mini: Open Model for Evaluating MCP Tool Use - Qualifire Blog, accessed on December 17, 2025, https://www.qualifire.ai/posts/ranger-mini-open-model-for-evaluating-mcp-tool-use
16. ADR Analysis MCP Server: A Deep Dive for AI Engineers, accessed on December 17, 2025, https://skywork.ai/skypage/en/adr-analysis-mcp-server-ai-engineers/1980183120013729792
17. Software Architecture Decision-Making Practices and Challenges: An Industrial Case Study, accessed on December 17, 2025, https://www.researchgate.net/publication/282272890_Software_Architecture_Decision-Making_Practices_and_Challenges_An_Industrial_Case_Study
18. What Is MCP (Model Context Protocol)? | Solo.io, accessed on December 17, 2025, https://www.solo.io/topics/ai-infrastructure/what-is-mcp
19. Model Context Protocol - Wikipedia, accessed on December 17, 2025, https://en.wikipedia.org/wiki/Model_Context_Protocol
20. Understanding MCP servers - Model Context Protocol, accessed on December 17, 2025, https://modelcontextprotocol.io/docs/learn/server-concepts
21. Beyond APIs: Lessons from Building with the Model Context Protocol (MCP) - Medium, accessed on December 17, 2025, https://medium.com/@wim.henderickx/beyond-apis-lessons-from-building-with-the-model-context-protocol-mcp-dc85ca83bb89
22. The Power of Model Context Protocol: Using Natural Language to Query GraphDB, accessed on December 17, 2025, https://graphwise.ai/blog/the-power-of-model-context-protocol-using-natural-language-to-query-graphdb/
23. Context-aware incident handling with MCP: A strategic view with a practical case, accessed on December 17, 2025, https://www.thoughtworks.com/en-us/insights/blog/generative-ai/context-aware-incident-handling-with-MCP-strategic-view-with-a-practical-case
24. Agentic AI: MCP & Its Impact on Observability Automation - Mezmo, accessed on December 17, 2025, https://www.mezmo.com/learn-observability/agentic-ai-what-is-model-context-

protocol-agent2agent-and-how-does-this-impact-automation

25. Understanding Cognitive Load in Software Engineering Teams and Systems - SoftwareSeni, accessed on December 17, 2025, https://www.softwareseni.com/understanding-cognitive-load-in-software-engineering-teams-and-systems/

26. Navigating Cognitive Load in Software Development: Experiences, Strategies, and Future Tools for Enhanced Mental Workflow - Zigpoll, accessed on December 17, 2025, https://www.zigpoll.com/content/how-do-software-developers-experience-and-manage-cognitive-load-during-complex-coding-tasks-and-what-tools-or-design-features-could-better-support-their-mental-workflow

27. How can cognitive load theory be applied to improve the usability of developer tools? Cognitive Load Theory (CLT) offers a powerful framework to optimize developer tool design by minimizing mental strain and maximizing productivity. Here's how applying CLT principles directly enhances developer tools' usability and developer experience (DX), helping teams build smarter, more effective - Zigpoll, accessed on December 17, 2025, https://www.zigpoll.com/content/how-can-cognitive-load-theory-be-applied-to-improve-the-usability-of-developer-tools

28. How Cognitive Load Affects Developer Productivity — and How to Fix It - iDelsoft, accessed on December 17, 2025, https://idelsoft.com/blog/tpost/how-cognitive-load-affects-developer-productivity-and-how-to-fix-it

29. SWE-Bench Pro: Can AI Agents Solve Long-Horizon Software Engineering Tasks? - arXiv, accessed on December 17, 2025, https://arxiv.org/html/2509.16941v1

30. True Cost of AI-Generated Code. A Strategic Analysis of "Comprehension... | by Justin Hamade | Medium, accessed on December 17, 2025, https://medium.com/@justhamade/true-cost-of-ai-generated-code-f4362391790c

31. Domain-Driven Design (DDD) - Redis, accessed on December 17, 2025, https://redis.io/glossary/domain-driven-design-ddd/

32. Domain Driven Design in AI-Driven Era - DEV Community, accessed on December 17, 2025, https://dev.to/aws-heroes/domain-driven-design-in-ai-driven-era-4l3h

33. What is Ubiquitous Language? - Dremio, accessed on December 17, 2025, https://www.dremio.com/wiki/ubiquitous-language/

34. Revolutionizing Enterprise AI: Applying Domain-Driven Design for Agentic Applications, accessed on December 17, 2025, https://sathiyan.medium.com/revolutionizing-enterprise-ai-applying-domain-driven-design-for-agentic-applications-aa321fb991f4

35. Applying domain-driven design principles to multi-agent AI systems - James Croft, accessed on December 17, 2025, https://www.jamescroft.co.uk/applying-domain-driven-design-principles-to-multi-agent-ai-systems/

36. The Rise of Agentic IDEs: What Cursor, Windsurf, and Others Tell Us About the Future of Development | by Anand Satheesh Kumar Nair | Medium, accessed on

December 17, 2025, https://medium.com/@anandnair/the-rise-of-agentic-ides-what-cursor-windsurf-and-others-tell-us-about-the-future-of-development-bb36b45e0701

37. This is a experimental code repository to try out various semantic code search techniques. - GitHub, accessed on December 17, 2025, https://github.com/chamikabm/semantic-code-search//

38. ChunkHound | Awesome MCP Servers, accessed on December 17, 2025, https://mcpservers.org/servers/ofriw/chunkhound

39. agentic-software-engineering-for-leaders-handbook/Agentic_Software_Engineering_for_Leaders_Handbook.md at main - GitHub, accessed on December 17, 2025, https://github.com/lucianoayres/agentic-software-engineering-for-leaders/blob/main/Agentic_Software_Engineering_for_Leaders_Handbook.md

40. Why AI Coding Tools Make Experienced Developers 19% Slower and How to Fix It, accessed on December 17, 2025, https://www.augmentcode.com/guides/why-ai-coding-tools-make-experienced-developers-19-slower-and-how-to-fix-it

41. 10 practical tips for AI-assisted development, accessed on December 17, 2025, https://www.civic.com/resources/tips-for-ai-assisted-development

42. (PDF) Using SOLO taxonomy to explore students' mental models of the programming variable and the assignment statement - ResearchGate, accessed on December 17, 2025, https://www.researchgate.net/publication/235920619_Using_SOLO_taxonomy_to_explore_students'_mental_models_of_the_programming_variable_and_the_assignment_statement

43. How cognitive complexity creates hidden friction in engineering organizations - DX, accessed on December 17, 2025, https://getdx.com/blog/cognitive-complexity/

44. Production-Ready AI Agents: 8 Patterns That Actually Work (with Real Examples from Bank of America, Coinbase & UiPath) | by Sai Kumar Yava | Nov, 2025 | Towards AI, accessed on December 17, 2025, https://pub.towardsai.net/production-ready-ai-agents-8-patterns-that-actually-work-with-real-examples-from-bank-of-america-12b7af5a9542

45. Customize agent workflows with advanced orchestration techniques using Strands Agents, accessed on December 17, 2025, https://aws.amazon.com/blogs/machine-learning/customize-agent-workflows-with-advanced-orchestration-techniques-using-strands-agents/

46. Top 10+ Agentic Orchestration Frameworks & Tools - Research AIMultiple, accessed on December 17, 2025, https://research.aimultiple.com/agentic-orchestration/

47. Diagnosing and Measuring AI Agent Failures: A Complete Guide - Maxim AI, accessed on December 17, 2025, https://www.getmaxim.ai/articles/diagnosing-and-measuring-ai-agent-failures-a-complete-guide/

48. SHIELDA: Structured Handling of Exceptions in LLM-Driven Agentic Workflows -

arXiv, accessed on December 17, 2025, https://arxiv.org/html/2508.07935v1

49. How Can Input Reformulation Improve Tool Usage Accuracy in a Complex Dynamic Environment? A Study on $\tau$-bench - arXiv, accessed on December 17, 2025, https://arxiv.org/html/2508.20931v2

50. Detecting Sleeper Agents in Large Language Models via Semantic Drift Analysis - arXiv, accessed on December 17, 2025, https://arxiv.org/abs/2511.15992

51. Drift Detection in Large Language Models: A Practical Guide | by Tony Siciliani | Medium, accessed on December 17, 2025, https://medium.com/@tsiciliani/drift-detection-in-large-language-models-a-practical-guide-3f54d783792c

52. How to Measure Drift in ML Embeddings - Towards Data Science, accessed on December 17, 2025, https://towardsdatascience.com/how-to-measure-drift-in-ml-embeddings-ee8adfe1e55e/

53. Fitness Functions for Your Architecture - InfoQ, accessed on December 17, 2025, https://www.infoq.com/articles/fitness-functions-architecture/

54. Agile Architecture Part.7: Fitness Functions & Architectural Guardrails | by Stefano Rossini, accessed on December 17, 2025, https://medium.com/@stefano.rossini.mail/agile-architecture-part-7-fitness-functions-architectural-guardrails-28ed22ade476

55. ArchUnit: Unit test your Java architecture, accessed on December 17, 2025, https://www.archunit.org/

56. How to Ensure Clean Architecture with ArchUnit? | by Akiner Alkan | Java and Beyond, accessed on December 17, 2025, https://medium.com/java-and-beyond/ensuring-clean-architecture-with-archunit-91d43959e648

57. Are LLMs Truly Solving Software Problems — or Are Agents Doing It? - SambaNova, accessed on December 17, 2025, https://sambanova.ai/blog/are-llms-truly-solving-software-problems

58. Scoring 71% on SWE-bench Verified in half the steps - Vibe Kanban, accessed on December 17, 2025, https://www.vibekanban.com/blog/scoring-71-on-swe-bench-verified-in-half-the-steps

59. An Empirical Study on Failures in Automated Issue Solving - arXiv, accessed on December 17, 2025, https://arxiv.org/html/2509.13941v1

60. What Is a Prompt Injection Attack? - IBM, accessed on December 17, 2025, https://www.ibm.com/think/topics/prompt-injection

61. Indirect Prompt Injection: The Hidden Threat Breaking Modern AI Systems - Lakera, accessed on December 17, 2025, https://www.lakera.ai/blog/indirect-prompt-injection

62. PromptPwnd: Prompt Injection Vulnerabilities in GitHub Actions Using AI Agents - Aikido, accessed on December 17, 2025, https://www.aikido.dev/blog/promptpwnd-github-actions-ai-agents

63. Output-Invariant and Time-Based Testing – Practical Techniques for Black-Box Enumeration of LLMs - TrebledJ's Pages, accessed on December 17, 2025,

https://trebledj.me/posts/output-invariant-prompt-injection/

64. OWASP Top 10 2025: How to Operationalize Software Supply Chain Security for Developer Environments - Safety's cybersecurity, accessed on December 17, 2025, https://www.getsafety.com/blog-posts/owasp-top-10-2025

65. Supply Chain Attacks Through Open Source Software: A Comprehensive Analysis of NPM, PyPI, and Docker Hub Vulnerabilities - ODU Digital Commons, accessed on December 17, 2025, https://digitalcommons.odu.edu/cgi/viewcontent.cgi?article=1143&context=covacci-undergraduateresearch

66. Understanding and Preventing Dependency Confusion Attacks - OX Security, accessed on December 17, 2025, https://www.ox.security/blog/understanding-and-preventing-dependency-confusion-attacks/

67. Prompt injection to RCE in AI agents - The Trail of Bits Blog, accessed on December 17, 2025, https://blog.trailofbits.com/2025/10/22/prompt-injection-to-rce-in-ai-agents/

68. Assessing architectural drift in commercial software development: A case study | Request PDF - ResearchGate, accessed on December 17, 2025, https://www.researchgate.net/publication/220280551_Assessing_architectural_drift_in_commercial_software_development_A_case_study

69. Principles and Strategies for AI-Augmented Human Reasoning - arXiv, accessed on December 17, 2025, https://arxiv.org/pdf/2503.15530