

Building the Librarian's Brain

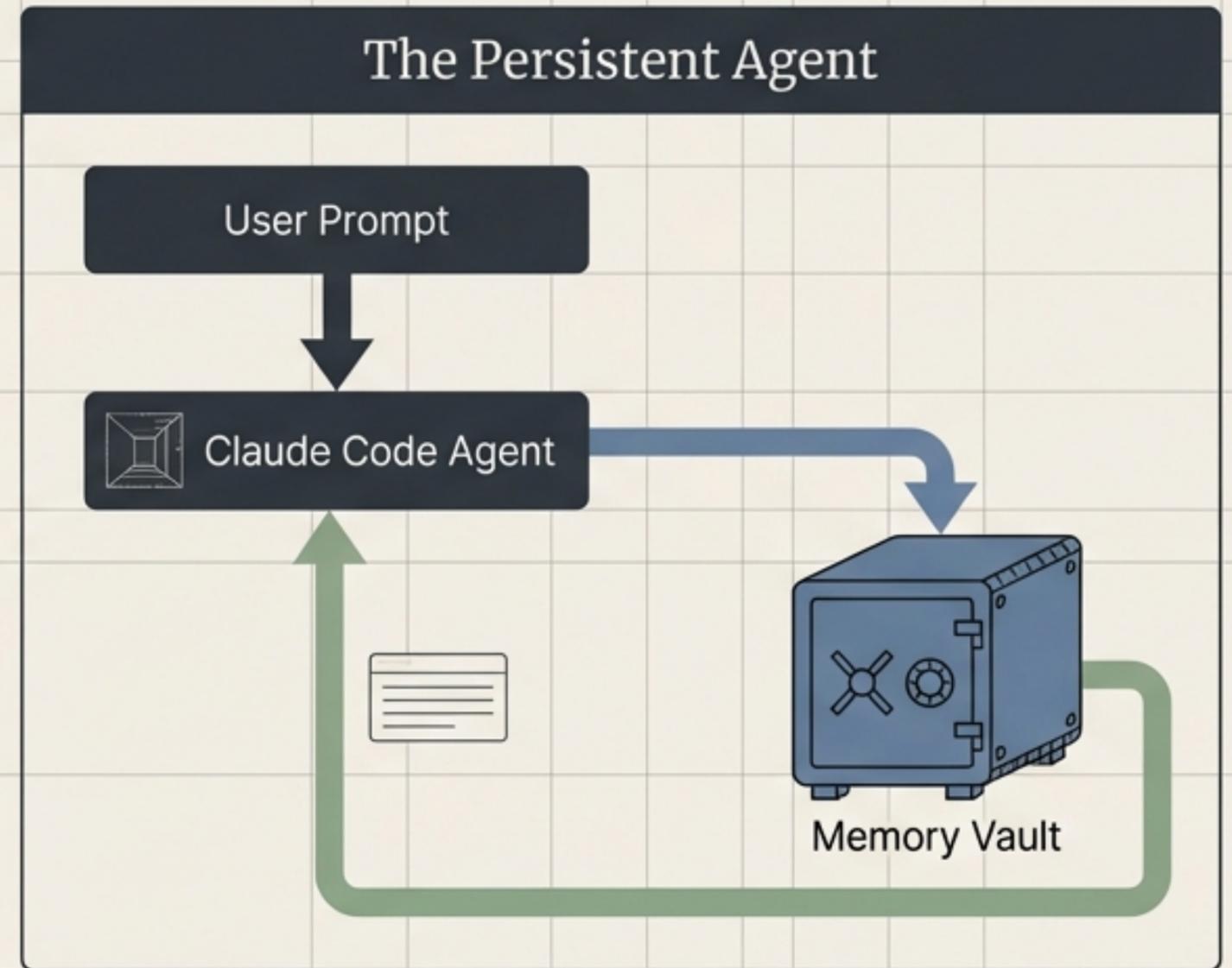
CC-303 Assessment Guide:
Persistent Memory & Context
for Autonomous Agents

Mastering cross-session
persistence, vector storage,
and semantic recall.

The Context Problem: Stateless vs. Stateful Automation



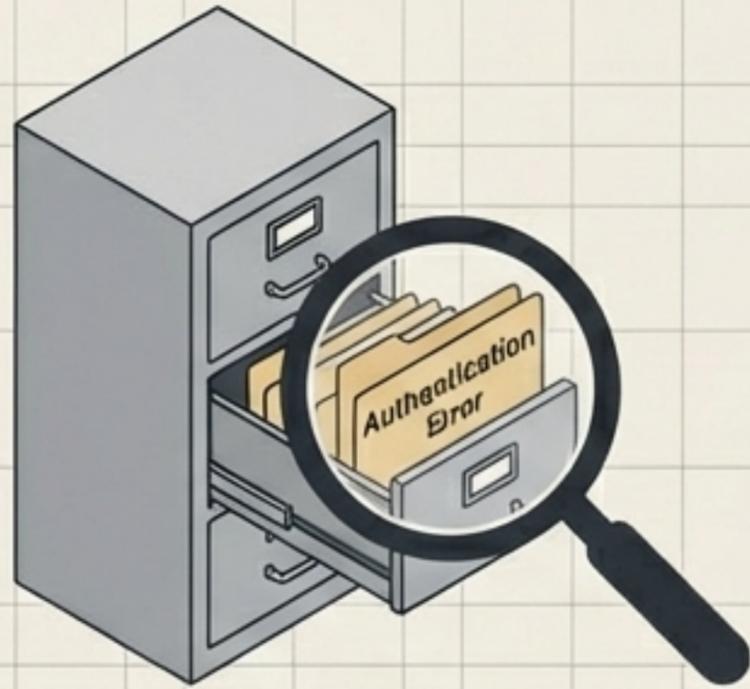
Every new session starts at zero. Architectural decisions, user preferences, and project history are lost.



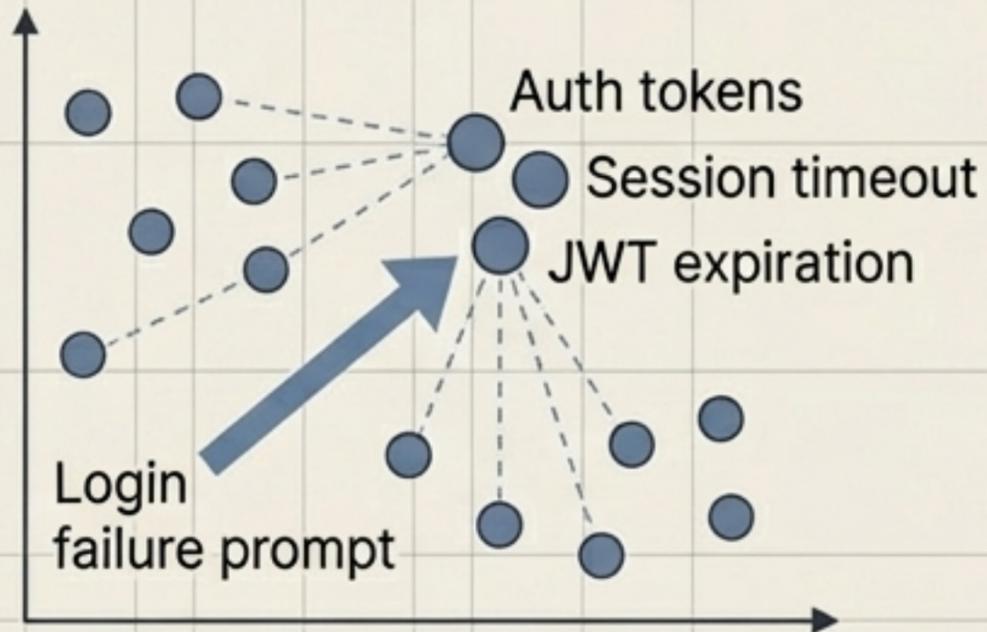
Context compounds infinitely. The agent learns the codebase permanently and auto-injects historical decisions into new sessions.

Without persistent memory, Claude Code is a temporary assistant. With it, it becomes a **continuous teammate**.

The Librarian Metaphor: Finding Meaning, Not Matches

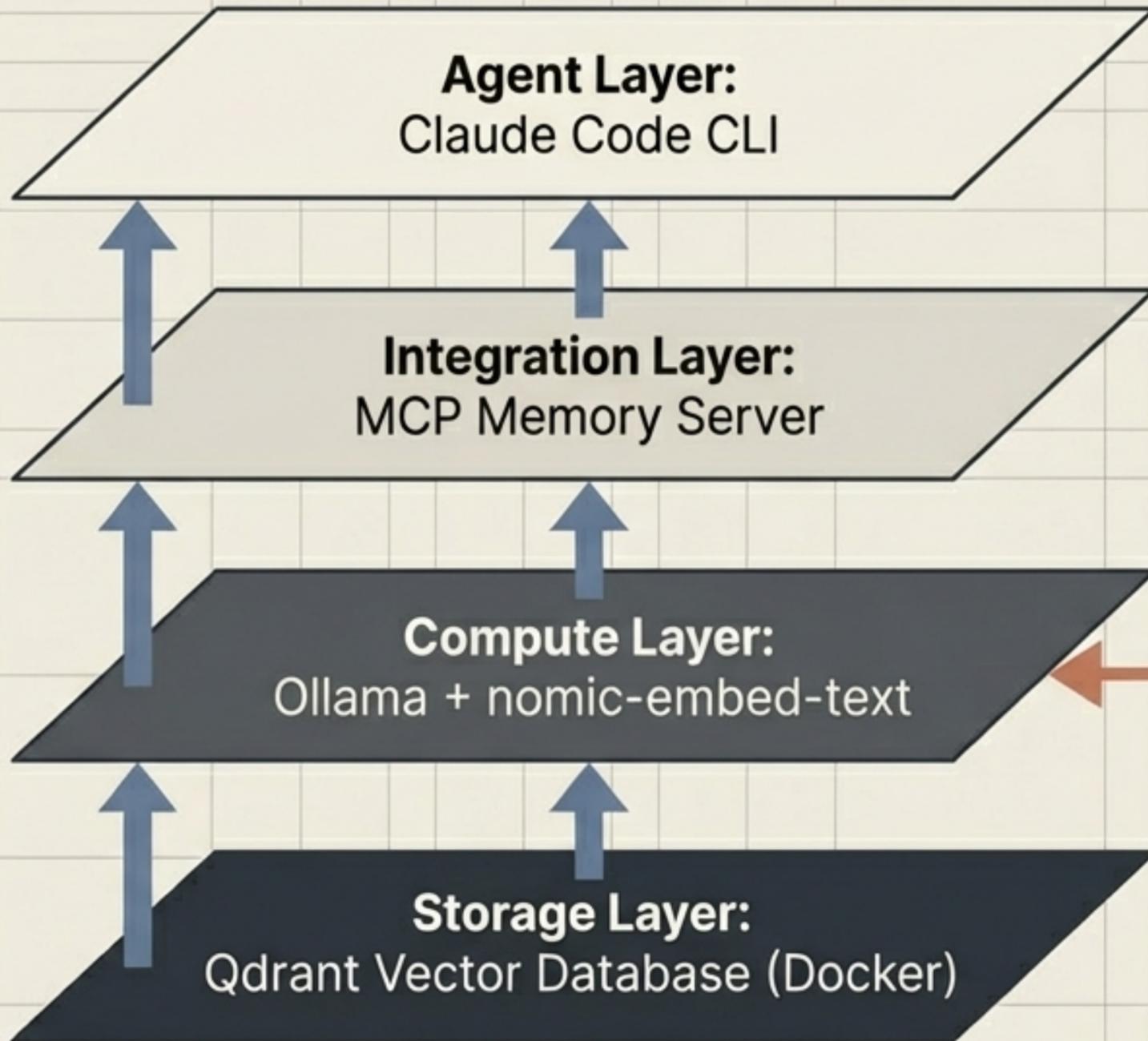


Traditional databases find records by exact string match. A search for 'caching layer' misses 'Redis implementation'.



Vector databases act like a master librarian. They map information in a high-dimensional space, retrieving memories based on conceptual proximity. Claude Code recalls context by meaning.

The CC-303 Memory Architecture Stack



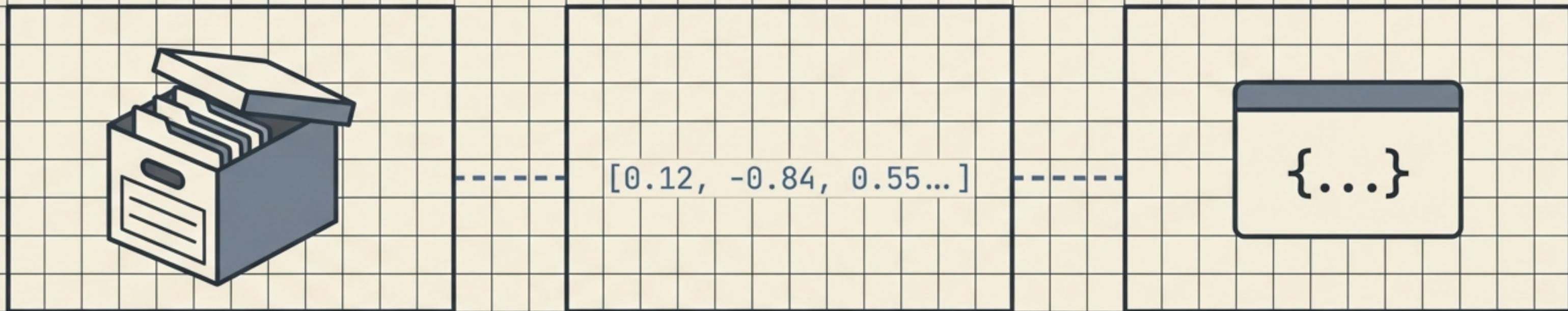
Local: All components run entirely on your machine.

Private: No memory data leaves your environment.

Persistent: Data survives process termination.

Crucial Lab Requirement: Embeddings are generated locally to preserve privacy and reduce latency.

Anatomy of the Vector Database (Qdrant)



Collections

The overall namespace. Memories are grouped into project-specific collections to prevent context cross-contamination.

Points

The individual vectors. High-dimensional arrays of floats representing the semantic meaning of the memory.

Payloads

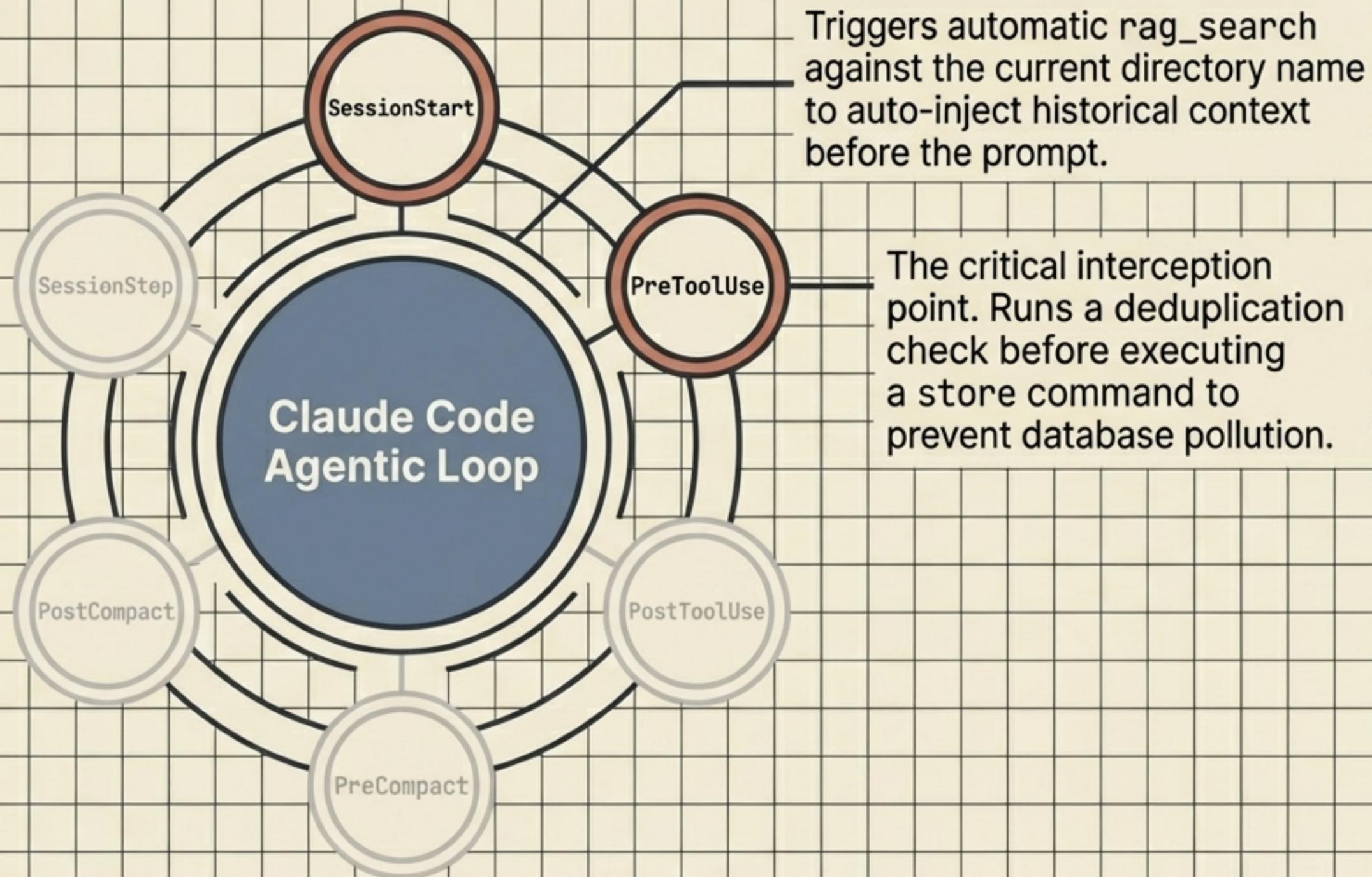
Attached metadata. Contains timestamps, original raw text, and the ID of the agent that generated the memory.

The 7 Tools of the MCP Memory Server

| Storage Operations | Retrieval Operations | Synthesis Operations |
|---|--|--|
| <p data-bbox="503 652 776 780"><code>store</code></p> <p data-bbox="243 846 1042 977">Saves a new concept or fact into the vector space.</p> <hr data-bbox="193 1046 1086 1050"/> | <p data-bbox="1502 652 1825 780"><code>recall</code></p> <p data-bbox="1269 846 2059 977">Exact-match fetching based on payload metadata.</p> <hr data-bbox="1219 1046 2112 1050"/> <p data-bbox="1419 1121 1909 1249"><code>rag_search</code></p> <p data-bbox="1269 1309 2059 1440">Semantic fetching based on conceptual proximity.</p> | <p data-bbox="2502 652 2868 780"><code>episode</code></p> <p data-bbox="2485 808 2885 936"><code>learning</code></p> <p data-bbox="2469 964 2902 1091"><code>procedure</code></p> <p data-bbox="2452 1119 2918 1247"><code>trajectory</code></p> <p data-bbox="2202 1309 3178 1515">Higher-order tools. Construct complex memory structures by extracting reusable patterns over time.</p> |

MCP bridges the gap between natural language prompts and Qdrant's REST API.

The Memory Hook Lifecycle (6-Event Flow)



Memory is not a passive database. It actively intercepts the agent's lifecycle to dynamically inject and protect context.

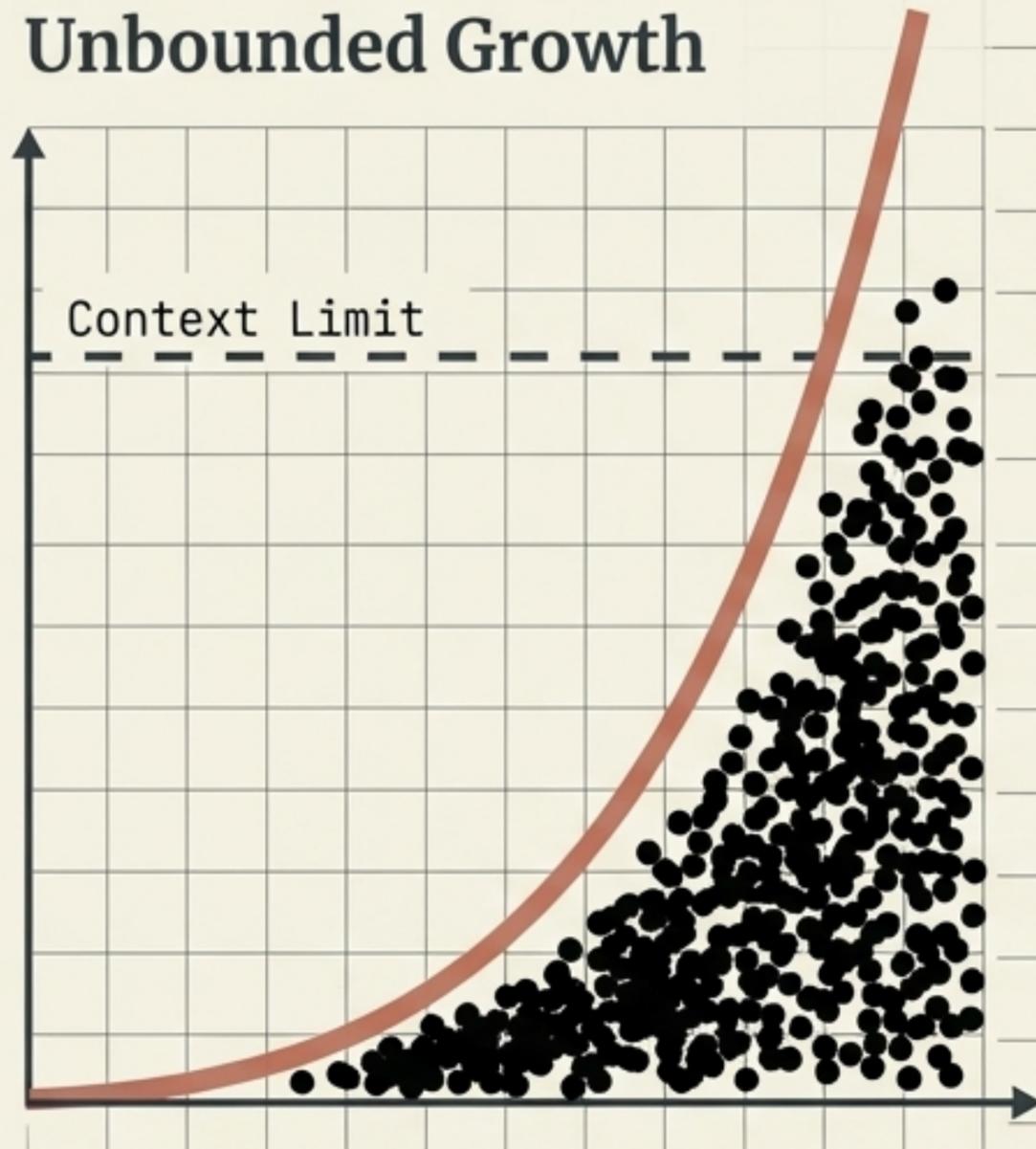
Diagnostic Gauge: Vector Similarity Thresholds



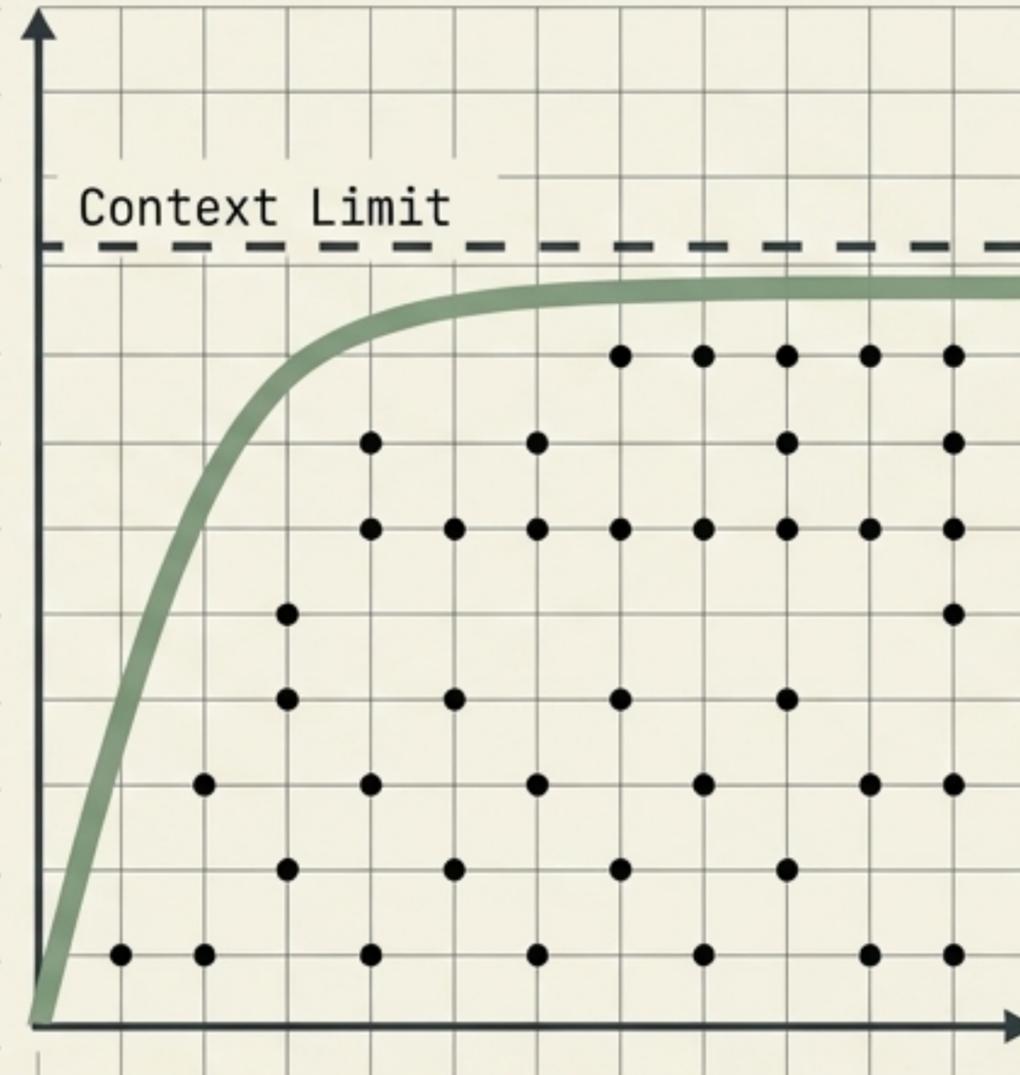
CC-303 Pitfall: Treating 0.80 as a duplicate will prematurely block writes.
The system must know when to link versus when to drop.

Deduplication & Pruning Strategies

Unbounded Growth



Curated Memory



If every store command succeeds, the vector database becomes polluted with identical entries, exhausting the LLM's context window.

The Solution

1. Agent invokes store.
2. PreToolUse hook intercepts.
3. Hook runs a rag_search for the exact concept.
4. If similarity > 0.92, hook blocks tool execution and returns existing memory ID.

Governance Integration: Fail-Open vs. Fail-Closed

Fail-Open

Scenario 1: Deduplication Error

Trigger: Qdrant is temporarily unreachable during a PreToolUse deduplication check.



Action: Proceed with Storage

Rationale: Better to accept potential duplicates than to permanently lose valuable user data.

Fail-Closed

Scenario 2: Classification Error

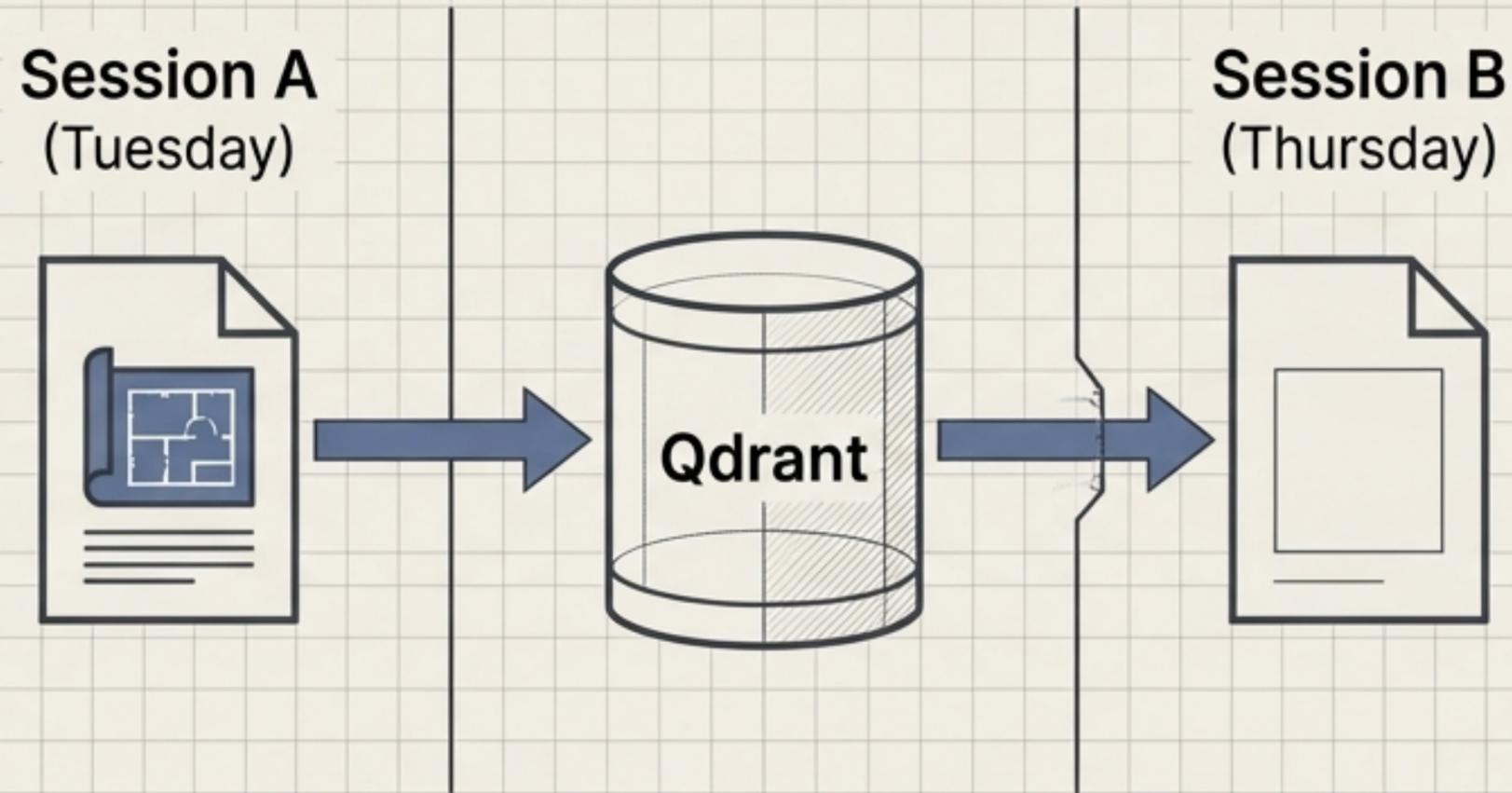
Trigger: The Memory Governor encounters a regex error while checking for restricted data.



Action: Halt & Block Storage

Rationale: Never risk storing Restricted content (e.g., API keys) in a Public collection. Unconditional block.

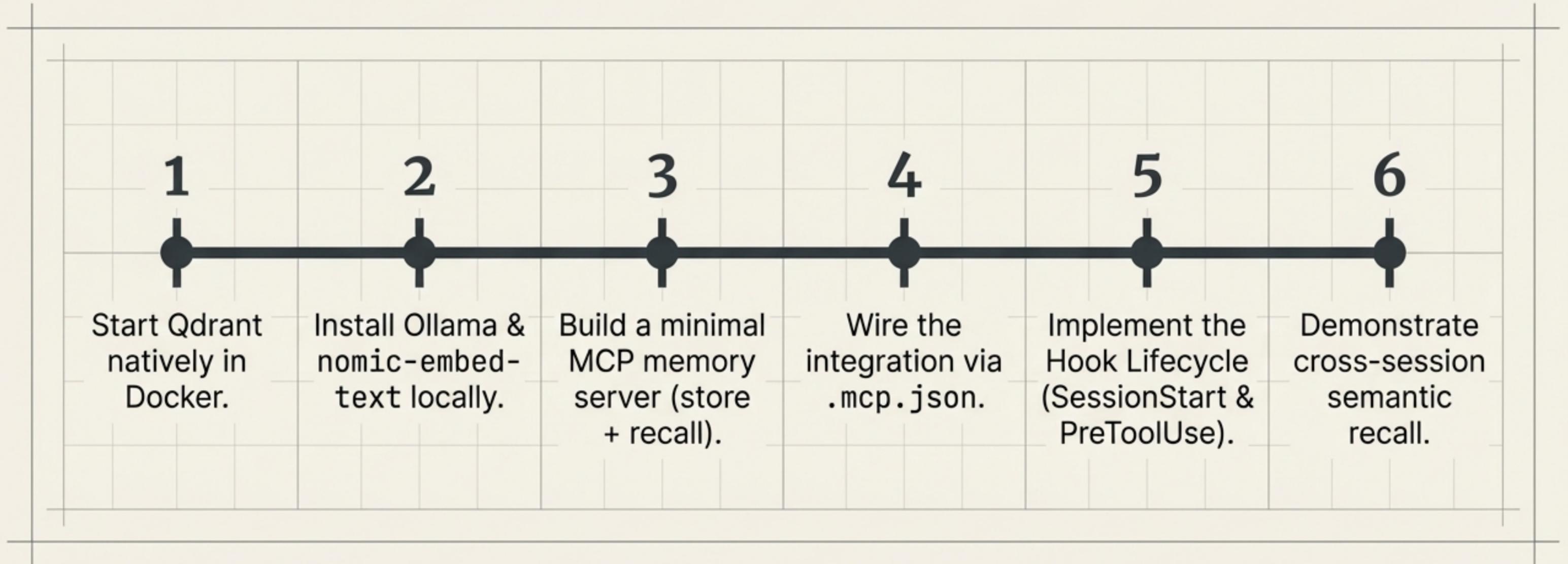
Cross-Session Persistence & Sync



Core Mechanism: Additive Merge

- Memory doesn't just span time; it syncs state.
- Launching `claude --continue` or starting a fresh session automatically pulls historical trajectory data.
- Result: An architectural decision made in Session A on **Tuesday** is automatically recalled in Session B on **Thursday** without the user restating the requirements.

CC-303 Practical Lab Progression



Top Assessment Pitfalls (Instructor Diagnostics)



Architectural Misstep

Error: Running Ollama inside Docker.

Correction: Causes massive performance degradation. Ollama must be run natively on the host OS.



Threshold Confusion

Error: Rejecting memories at > 0.80 similarity.

Correction: 0.80 is the linking threshold. Only reject exact duplicates at > 0.92 .

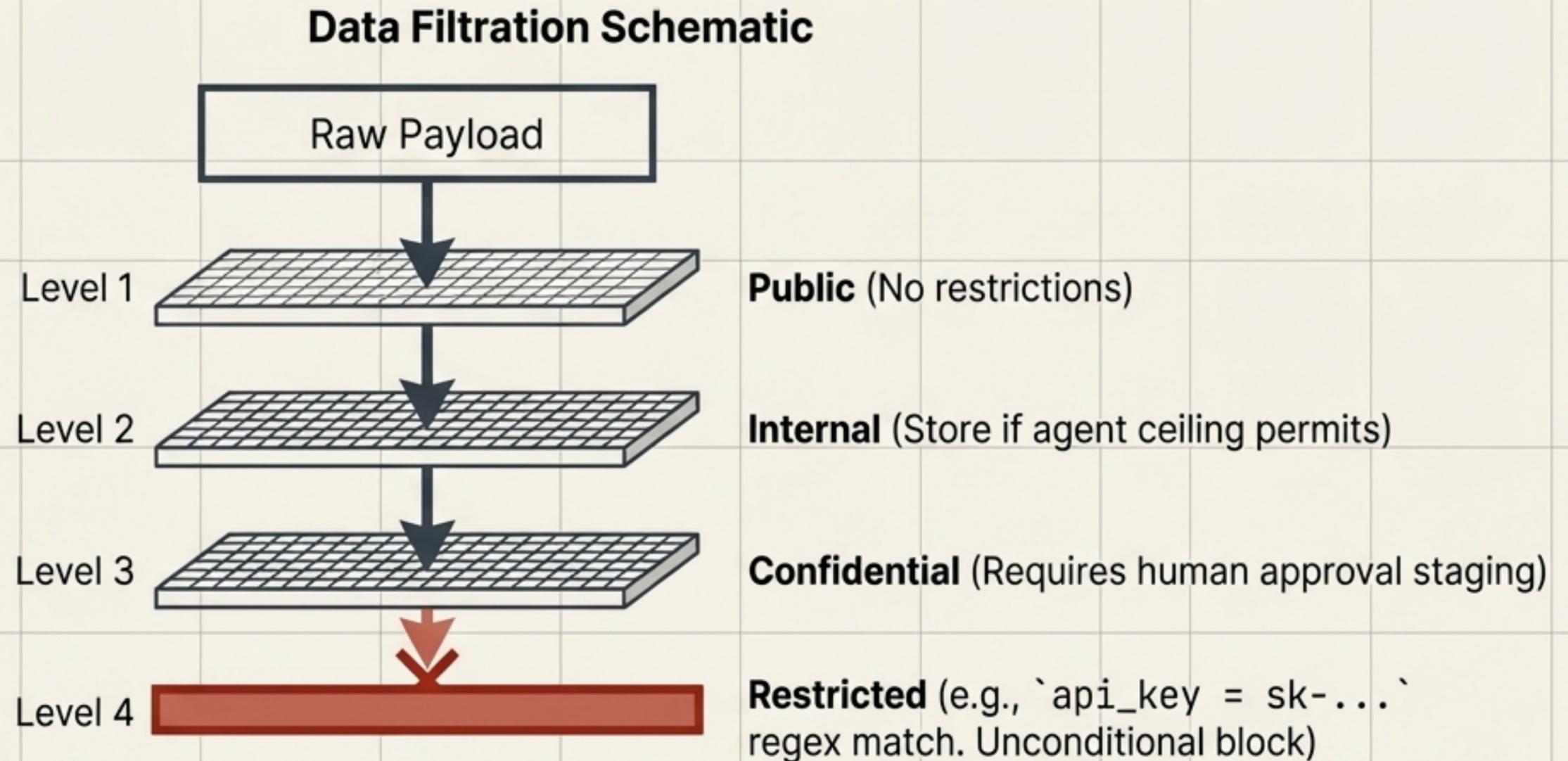


Unbounded Collections

Error: Failing to implement the deduplication interception.

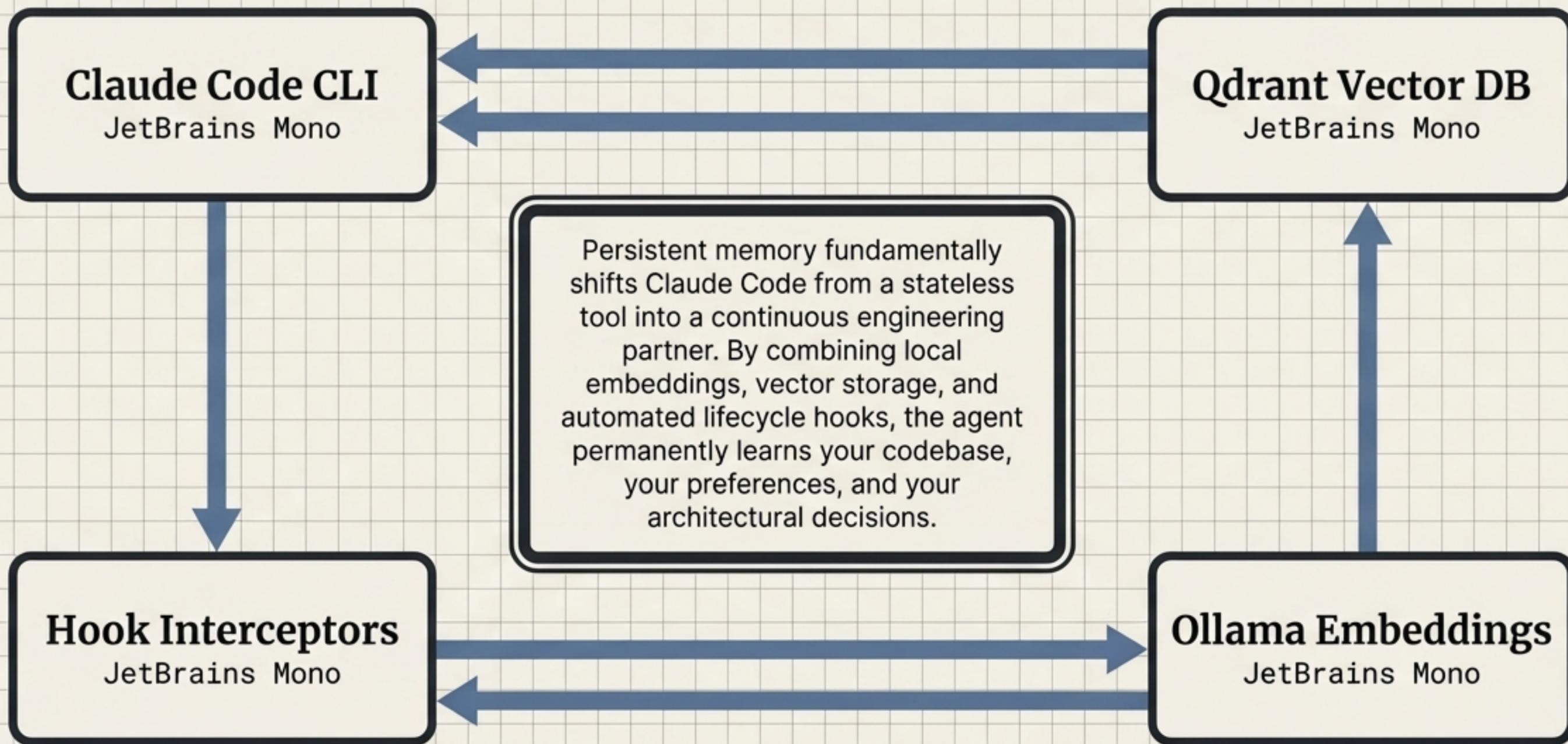
Correction: A PreToolUse hook must be active to catch and prune duplicate store commands before they hit the database.

Integrating the Memory Governor



Takeaway: The Memory Governor scans all payloads prior to storage. Even a Trust Level 5 agent cannot write **Restricted** content to the vector database.

Synthesis: The Infinite Agent



Memory is the bridge between automation and autonomy.