# Unified Agentic Memory Layer for AI Coding Agents

*Enterprise AI Development Infrastructure + Research■Grade Agent Framework*

Version 1.0 • Generated 2026-01-06

# Executive Summary

Modern coding agents (Cursor, Claude/Claude Code, Codex, Kiro/Kiro CLI, and future IDE/CLI/Cloud agents) are converging on a shared problem: agents can act, but they cannot reliably share durable project knowledge, decisions, and coordination state across tools, sessions, and parallel workstreams. This paper proposes a vendor■agnostic "agentic unified memory layer" that provides persistent code intelligence, structured decision logs, and interoperable connectivity so that any coding agent can plug in and immediately become aware of what the others have done—while preserving auditability and governance.

The design is intentionally protocol■centric: we treat agents as interchangeable clients and standardize on open interoperability primitives (A2A, MCP, Agent Protocol, plus conventional HTTP/gRPC/event streams) to maximize ecosystem compatibility and enterprise adoption.

# 1. Problem Statement

AI coding assistants today remain fragmented and frequently stateless. Without a shared memory substrate, parallel agents duplicate work, introduce contradictory changes, and lose the rationale behind architectural decisions. Enterprises also lack consistent governance controls (who changed what, why, and under which policy).

## 1.1 Goals

**G1.** Any agent (IDE, CLI, cloud sandbox) can **read** project context and **write** decisions and changes.
**G2.** Support **semantic recall** (vector) plus **structured reasoning** (graph) with time■aware audit trails.
**G3.** Provide **multi■agent coordination** (locks, leases, plans) to safely run agents in parallel.
**G4.** Enterprise-grade security: authN/authZ, tenancy, provenance, and policy enforcement.
**G5.** Research-grade observability: reproducible runs, reasoning summaries, dataset export, evaluation hooks.

# 2. Reference Architecture

The platform consists of three planes: (A) Ingestion & Indexing (code $\rightarrow$ symbols/graphs/vectors), (B) Memory & Knowledge (graph + vector + time), and (C) Connectivity & Coordination (protocol adapters, agent runtime contracts, eventing, and governance).

## 2.1 System Diagram (Mermaid)

```
flowchart TB
  subgraph Connect["Connectivity & Coordination"]
    AP["Agent Protocol API\n(runs/threads/store)"]
    A2A["A2A Server\n(agent.json + JSON-RPC/SSE)"]
    MCP["MCP Server/Client\n(tools/resources)"]
    BUS["Event Bus\n(Kafka/NATS/RabbitMQ)"]
```

```
    end

    subgraph Mem["Unified Memory & Knowledge"]
      VEC["Vector Index\n(embeddings + filters)"]
      GRA["Knowledge Graph\n(symbols/decisions/edges)"]
      AUD["Temporal Audit Log\n(changesets, provenance)"]
    end

    subgraph Ingest["Ingestion & Indexing"]
      TS["Tree-sitter/AST Parsers"]
      DIFF["Git Diff + Blame + PR events"]
      DOC["Docs/Issues/ADRs/AGENTS.md"]
    end

    Agents["Coding Agents\n(Cursor, Claude, Codex, Kiro, etc.)"] --> Connect
    Connect --> Mem
    Ingest --> Mem
```

# 3. Agent Connectivity: Protocols and Adapters

Your platform should not bet on a single agent vendor or IDE. The safest strategy is to implement a small set of well-chosen interoperability surfaces and provide adapters per agent environment (IDE extensions, CLI wrappers, cloud sandboxes). Below are the primary options, organized by how "agent-native" they are.

## 3.1 Agent-Native Interoperability Standards

**Agent2Agent (A2A)**: an open protocol enabling agents on different stacks to discover capabilities (via an Agent Card at */.well-known/agent.json*), exchange tasks, and stream updates—built on familiar web standards such as JSON-RPC and SSE. A2A is specifically positioned for cross-platform, multi-agent collaboration.

**Model Context Protocol (MCP)**: an open standard for secure, two-way connections between AI tools and external data sources/tools via MCP servers and clients. For unified memory, MCP is ideal for exposing "memory tools" (store, search, retrieve, write-decision) in a standardized way, letting any MCP-capable agent consume the same memory substrate.

**LangChain Agent Protocol**: a production API contract centered on *Runs*, *Threads*, and *Store* (long-term memory) endpoints. This is a strong choice when you want framework-agnostic serving for agents and a standardized long-term store API.

**AGENTS.md**: a simple, open convention for placing agent guidance in repositories. Treat AGENTS.md as a first-class ingestion target and propagate its policies into the memory layer so every agent gets consistent project instructions.

## 3.2 Conventional Protocol Surfaces

**HTTP/REST** (lowest friction), **gRPC** (typed + fast), **WebSockets/SSE** (streaming), and **Message Bus** (pub/sub for events). These are useful to implement once as the internal transport layer even if you expose agent-native protocols externally.

## 3.3 Protocol Comparison Matrix

| Protocol | Best For | Strengths | Risks / Gaps |
|---|---|---|---|
| A2A | Agent-to-agent collaboration | Discovery via Agent Card; task lifecycle; streaming updates | Ecosystem still maturing; needs good authZ story |

| MCP | Tools/resources/context access | Standard tool schemas; secure 2-way connections; broad adoption | Not a full orchestration protocol by itself |
|---|---|---|---|
| Agent Protocol | Serving agents + runs/threads | Clean production API contract; long-term store | Doesn't define tool ecosystem like MCP |
| HTTP/REST | Universal integration | Everyone supports it; easiest debugging | No built-in semantics; must define schemas |
| gRPC | High-throughput internal RPC | Typed; efficient; multi-language | Harder to debug; requires codegen |
| Event Bus | Coordination & observability | Decoupled; scalable; replayable streams | Operational overhead; ordering/consistency comple |

# 4. Unified Memory: What Must Be Stored

A useful unified memory layer stores more than embeddings of code. It must capture **intent**, **decisions**, **provenance**, and **coordination state** so multiple agents can safely build on each other's work.
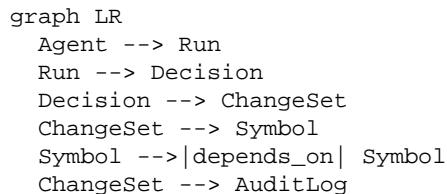
## 4.1 Memory Objects (Recommended)

• **Symbol Memory**: file/module/function/class nodes; call/import edges; ownership; version hash.
• **Decision Memory**: ADR-like records: problem, options, rationale, trade-offs, chosen path.
• **ChangeSets**: diffs grouped with a decision; includes tests run, CI results, PR links, rollbacks.
• **Agent Runs**: prompts/tools used (summarized), retrieved context IDs, outcomes, confidence.
• **Policies**: AGENTS.md + org rules + security constraints; enforce at retrieval time.
• **Coordination Artifacts**: locks/leases, task claims, plans, checkpoints, unresolved conflicts.

## 4.2 Hybrid Retrieval Strategy

**Vector retrieval** provides semantic recall ("find related bug fixes"), while **graph traversal** provides structural reasoning ("show all symbols impacted by this module"), and **temporal queries** provide auditability ("what changed since last release and why").

### *Mermaid – Decision + ChangeSet Linking*

```
graph LR
  Agent --> Run
  Run --> Decision
  Decision --> ChangeSet
  ChangeSet --> Symbol
  Symbol -->|depends_on| Symbol
  ChangeSet --> AuditLog
```

# 5. Memory Frameworks and Boosters

You can implement memory primitives yourself, but frameworks can accelerate correctness and adoption. The best approach is to define a stable internal memory API and allow plug-in backends (Mem0, custom, etc.).

## 5.1 Mem0 (Recommended Core Memory Abstraction)

Mem0 is a long-term memory framework for agents. Research results report substantially lower p95 latency and large token savings compared to naive "stuff the whole chat" approaches, while improving long-dialog coherence. Architecturally, Mem0 supports multiple memory types and can sit above your vector/graph stores as the memory orchestration layer.

## 5.2 Orchestration Framework Options

• **LangGraph/LangChain**: strong graph-based orchestration; production patterns around persistence and threads.
• **AutoGen**: flexible multi-agent conversations and role-based collaboration patterns.
• **CrewAI**: lightweight multi-agent teams and workflows; easy "crew" constructs.
• **Pydantic-AI**: robust structured outputs + tool schemas (excellent for deterministic memory writes).

# 6. Multi-Agent Coordination at Scale

Parallel agents require explicit coordination primitives—otherwise they conflict, duplicate work, or converge on inconsistent architecture. The unified layer should provide **soft locks**, **leases**, **task claims**, and **conflict records** backed by the knowledge graph.

## 6.1 Coordination Models

**Pull**: agent queries memory when needed (low overhead).
**Push**: agent writes decisions/changes immediately (strong audit).
**Pub/Sub**: agents subscribe to events ("Symbol X changed", "Lock acquired"). Best for large teams and observability.

### *Mermaid – Lock Lease Workflow*

```
sequenceDiagram
  participant A as Agent A
  participant M as Memory Layer
  participant B as Agent B

  A->>M: Acquire lease(lock:Symbol Foo, ttl=5m)
  M-->>A: Lease granted (token)
  B->>M: Acquire lease(lock:Symbol Foo)
  M-->>B: Denied (active lease)
  A->>M: Commit ChangeSet + Decision + Release lease
  M-->>B: Event: lease released
```

# 7. Enterprise Controls: Security, Governance, Compliance

Enterprises need the unified memory to be safe by default. Implement **authN** (OIDC), **authZ** (RBAC/ABAC), **tenancy**, **encryption**, and **policy enforcement** at retrieval and write time. Treat memory writes like code changes: require provenance, signatures, and traceability.

## 7.1 Provenance and Audit Trail

Every write should include: agent identity, model version, tool calls (summarized), retrieved context IDs, and human approvals (if any). This makes agent activity reviewable and supports rollback and compliance reporting.

# 8. Implementation Roadmap (Recommended)

**Phase 0** — Spec: define canonical object model (Symbols, Decisions, ChangeSets, Runs, Policies).
**Phase 1** — Indexing: Tree-sitter ingestion + embeddings + basic vector retrieval.
**Phase 2** — Protocols: implement Agent Protocol (runs/threads/store) + MCP server for memory tools + optional A2A gateway.
**Phase 3** — Graph: add symbol graph + decision edges + temporal audit queries.
**Phase 4** — Coordination: leases/locks + event bus + conflict resolution UI.
**Phase 5** — Enterprise: tenancy, policy engine, observability, evaluation harness.

# 9. Best-Overall Approach (Scaled Recommendation)

If you want the best blend of ecosystem compatibility and enterprise reliability, build a **protocol-first memory service** with three public surfaces:
1) **MCP** for standardized tools/resources access, 2) **Agent Protocol** for runs/threads/store lifecycle, and 3) optional **A2A** gateway for agent-to-agent delegation and discovery. Internally, use a hybrid vector+graph design with a time-aware audit log, and optionally leverage Mem0 as the memory orchestration layer on top of your stores.

# References (Selected)

• **Agent2Agent (A2A) announcement** — Google Developers Blog.
https://developers.googleblog.com/en/a2a-a-new-era-of-agent-interoperability/

• **A2A codelab (agent.json discovery)** — Google Codelabs.
https://codelabs.developers.google.com/intro-a2a-purchasing-concierge

• **Model Context Protocol announcement** — Anthropic. https://www.anthropic.com/news/model-context-protocol

• **MCP specification (authoritative)** — modelcontextprotocol.io.
https://modelcontextprotocol.io/specification/2025-11-25

• **LangChain Agent Protocol repo/docs** — LangChain. https://github.com/langchain-ai/agent-protocol

• **Mem0 paper (arXiv:2504.19413)** — arXiv. https://arxiv.org/abs/2504.19413

• **OpenAI introduces Codex** — OpenAI. https://openai.com/index/introducing-codex/

• **Codex CLI docs** — OpenAI Developers. https://developers.openai.com/codex/cli/

• **OpenAI + Linux Foundation AAIF** — OpenAI. https://openai.com/index/agentic-ai-foundation/

• **AGENTS.md standard** — agents.md. https://agents.md/

• **Cursor Codebase Indexing docs** — Cursor. https://cursor.com/docs/context/codebase-indexing

• **Kiro autonomous agent blog** — Kiro. https://kiro.dev/blog/introducing-kiro-autonomous-agent/

• **Anthropic donation of MCP + AAIF** — Anthropic.
https://www.anthropic.com/news/donating-the-model-context-protocol-and-establishing-of-the-agentic-ai-foundation