


AgenticVision-MCP: Persistent Visual Memory for AI Agents via the Model Context Protocol

Omoshola Owolabi 

Independent Researcher – AI/ML

owolabi.omoshola@outlook.com

February 19, 2026

Abstract

Large language models augmented with vision capabilities can perceive images within a single context window but lack persistent visual memory across sessions. When the context resets, all visual observations are lost. We present **AgenticVision-MCP**, a system that provides persistent, queryable visual memory to any LLM through the Model Context Protocol (MCP). The system captures images, generates CLIP embeddings, produces compact thumbnails, and stores all observations in a binary `.avis` file format with a 64-byte fixed header and JSON payload. We expose 10 tools, 6 resources, and 4 prompts over MCP’s JSON-RPC 2.0 transport, enabling any compliant client to capture screenshots, query visual history, compute visual diffs, find similar images by embedding cosine similarity, and link visual observations to a cognitive memory graph (AgenticMemory). In testing on an Apple M4 (macOS 26.2, Rust 1.90), the system passes 88 tests across 8 validation phases with zero failures: 35 Rust unit and edge-case tests, 47 Python client tests, and 6 cross-language integration tests including multi-agent scenarios where independent agents share a single `.avis` file. Per-capture storage overhead is 4.3 KB including a 512-dimensional embedding and JPEG thumbnail, enabling approximately 230,000 observations per gigabyte. All query operations complete within 1–2 ms on a 100-capture store (excluding process startup). The system requires zero external services and zero network dependencies.

1 Introduction

Vision-language models such as GPT-4o, Claude, and Gemini can perceive images within their context window, but this perception is ephemeral. When a conversation ends or context resets, all visual observations vanish. An agent that captured a screenshot thirty seconds ago cannot reference it after the session boundary. This is the *visual memory problem*: LLM-based agents perceive but do not remember what they see.

Current approaches to agent memory focus on text. Vector databases store text chunks as embeddings for retrieval-

augmented generation (RAG) [6]. Key-value stores like Mem0 [7] extract factual snippets. MemGPT [13] pages text in and out of context. None of these systems handle visual observations natively. Images are either discarded after the context window, or converted to text descriptions that lose the perceptual signal entirely.

The key insight behind AgenticVision-MCP is that *visual memory is a structured perception problem, not a text search problem*. An agent needs to store what it saw (the image), understand what it saw (the CLIP embedding), find what it saw before (cosine similarity search), detect what changed (pixel-level diff), and connect what it saw to what it knows (links to a cognitive memory graph). These operations require a purpose-built visual memory system.

We design AgenticVision-MCP as a standalone binary that communicates with any LLM client via the Model Context Protocol (MCP) [1]. The system consists of two Rust crates: a core vision library (`agentic-vision`) handling image processing, CLIP embedding, similarity search, and binary file I/O; and an MCP server (`agentic-vision-mcp`) exposing this functionality as 10 tools, 6 resources, and 4 prompts over stdio JSON-RPC transport. Visual observations are persisted in a single `.avis` binary file that requires no external database, no cloud service, and no network connection.

We validate the system through 8 phases of testing, achieving 88 passing tests with zero failures. We demonstrate multi-agent scenarios where independent agent processes share a single vision file, cross-system linking where visual captures connect to AgenticMemory [10] cognitive graph nodes, and sub-2ms query latency on real hardware. The remainder of this paper describes the architecture (Section 3), the MCP interface design (Section 4), evaluation results (Section 5), and discusses implications and limitations (Section 6).

2 Background and Related Work

Vision-Language Models. CLIP [15] demonstrated that contrastive pre-training produces visual embeddings that align with natural language, enabling zero-shot image

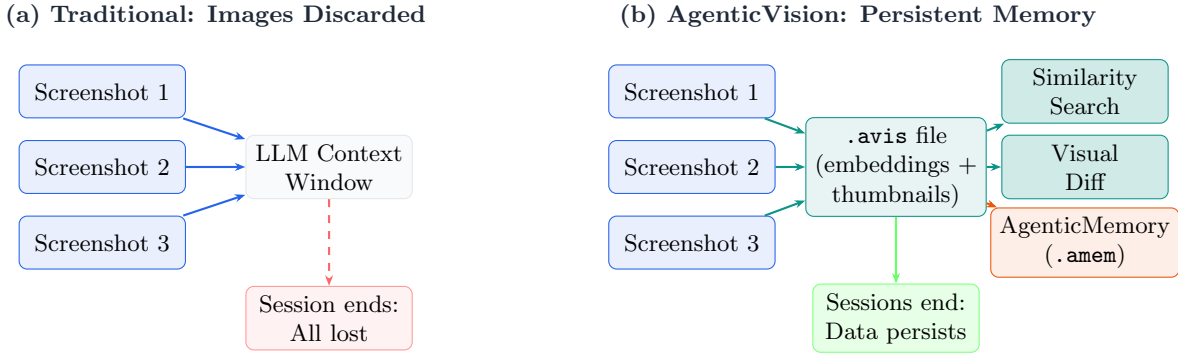


Figure 1: Motivating comparison. (a) Traditional LLMs perceive images within a context window but lose all visual data when the session ends. (b) AgenticVision-MCP persists every observation with its CLIP embedding and thumbnail in a binary `.avis` file, enabling similarity search, visual diff, and cross-session continuity. Observations can link to AgenticMemory cognitive graph nodes.

classification and cross-modal retrieval. Vision Transformers (ViT) [2] provide the backbone for modern visual encoders. AgenticVision-MCP uses CLIP ViT-B/32 to generate 512-dimensional embeddings per observation, enabling semantic similarity search.

Agent Memory Systems. Mem0 [7] provides key-value fact extraction from conversations. It excels at simple factual recall but offers no visual storage, no relationship tracking, and requires a cloud service. MemGPT [13] introduces virtual memory paging for LLM context windows, an elegant solution for text but inapplicable to visual data. OpenAI’s built-in memory [9] stores preference snippets but has no visual persistence. AgenticMemory [10] provides a binary cognitive graph with typed events and edges, designed for textual knowledge. None of these systems address visual perception.

Vector Databases. Pinecone [14], Weaviate [16], and FAISS [5] provide embedding storage and similarity search. While these could store image embeddings, they introduce significant operational complexity: cloud dependencies, API rate limits, and cost. They also lack the semantic structure needed for visual memory—there is no concept of “visual diff,” “capture session,” or “memory linking” in a generic vector store.

The Model Context Protocol. MCP [1] defines a standard interface between LLM clients and external tool servers. It specifies three primitives: *tools* (functions the LLM can invoke), *resources* (data the LLM can read), and *prompts* (guided workflows). MCP uses JSON-RPC 2.0 over stdio or HTTP+SSE transport. By implementing an MCP server, AgenticVision-MCP becomes accessible to any compliant client—Claude Desktop, VS Code, Cursor, or custom agents—without client-specific integration code.

Binary Formats. Protocol Buffers [3] and FlatBuffers [4] provide efficient binary serialization. The

Table 1: Comparison of agent memory and vision storage approaches.

System	Visual	Embed.	Diff	Deps	MCP
Mem0	×	×	×	Cloud	×
MemGPT	×	×	×	Python	×
OpenAI Mem.	×	×	×	Cloud	×
Pinecone	Partial		×	Cloud	×
FAISS	Partial		×	C++	×
AgenticMem.	×		×	None	
Ours				None	

`.avis` format draws inspiration from these systems but is purpose-built for visual memory: a fixed 64-byte header enables $O(1)$ metadata access, and the payload contains all observation data in a single self-contained file.

Table 1 compares AgenticVision-MCP with existing approaches.

3 Architecture

AgenticVision-MCP is organized as two Rust crates within a Cargo workspace. The core library (`agentic-vision`) handles image processing, embedding generation, similarity search, visual diff, and binary file I/O. The MCP server (`agentic-vision-mcp`) wraps this library with JSON-RPC 2.0 transport, tool dispatch, resource serving, and prompt expansion.

3.1 Visual Observation: The Atom

The fundamental data unit is a `VisualObservation`—a captured image enriched with computed features:

Listing 1: `VisualObservation` structure (simplified from Rust).

```
struct VisualObservation {
    id: u64,           // Unique capture ID
    timestamp: u64,    // Unix timestamp
```

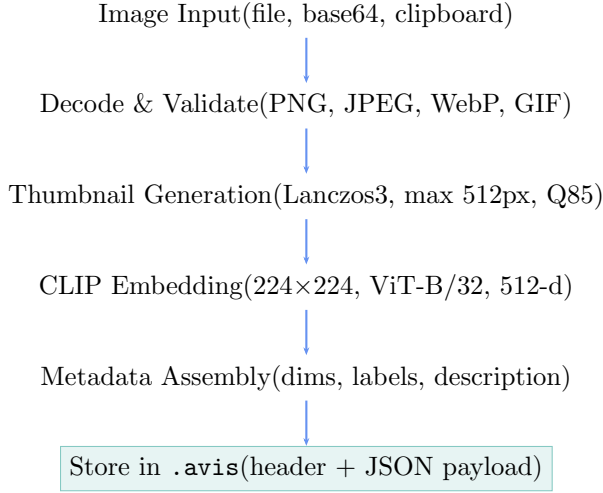


Figure 2: Capture pipeline. Each image passes through decoding, thumbnail generation (JPEG, Lanczos3 resampling), CLIP embedding (512-d vector via ONNX Runtime), metadata assembly, and storage in the binary `.avis` file.

```

{
  session_id: u32,    // Session scope
  source: CaptureSource, // File|Base64|
                      // Screenshot|Clipboard
  embedding: Vec<f32>, // 512-dim CLIP
  thumbnail: Vec<u8>, // JPEG, max 512x512
  metadata: ObservationMeta,
  memory_link: Option<u64>, // -> .amem node
}

```

The `embedding` field stores a 512-dimensional CLIP ViT-B/32 vector, L2-normalized after inference. The `thumbnail` field stores a JPEG-compressed image at quality 85, with the largest dimension capped at 512 pixels (aspect ratio preserved). The `memory_link` field enables cross-system integration by pointing to a node ID in an AgenticMemory `.amem` file.

3.2 The .avis Binary File Format

All visual observations are stored in a single binary file with the extension `.avis`. The format consists of a fixed 64-byte header followed by a JSON-serialized payload (Figure 3).

Design rationale. The fixed header enables fast metadata reads without parsing the payload. The JSON payload simplifies debugging and forward compatibility—new fields can be added without breaking readers. The single-file design means an agent’s visual memory is a single portable artifact: copy the `.avis` file to transfer all visual knowledge.

3.3 CLIP Embedding Engine

We use CLIP ViT-B/32 [15] via ONNX Runtime [8] to generate 512-dimensional visual embeddings. Input images are resized to 224×224 using Lanczos3 filtering,

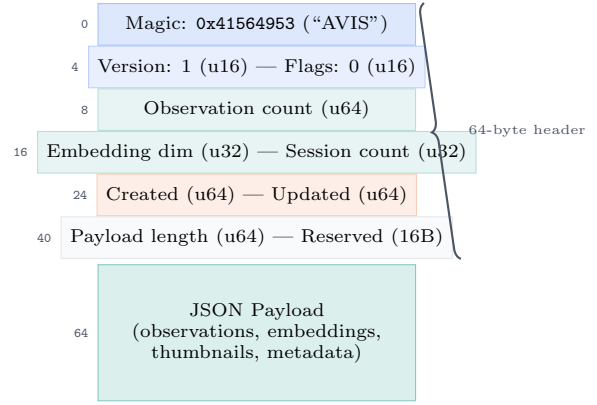


Figure 3: The `.avis` binary file format. A fixed 64-byte little-endian header provides $O(1)$ access to metadata (observation count, embedding dimension, timestamps). The payload contains JSON-serialized observations including embeddings and base64-encoded JPEG thumbnails.

converted to NCHW tensor layout, and normalized with CLIP’s published constants ($\mu = [0.481, 0.458, 0.408]$, $\sigma = [0.269, 0.261, 0.276]$). The output vector is L2-normalized. When the ONNX model is not available, the engine falls back to zero embeddings, allowing the system to function (without similarity search) on minimal installations.

3.4 Similarity Search

Given a query observation, we compute cosine similarity against all stored embeddings:

$$\text{sim}(\mathbf{a}, \mathbf{b}) = \frac{\sum_{i=1}^{512} a_i \cdot b_i}{\sqrt{\sum_{i=1}^{512} a_i^2} \cdot \sqrt{\sum_{i=1}^{512} b_i^2}} \quad (1)$$

Intermediate calculations use `f64` precision to avoid accumulation error over 512 dimensions; the final result is returned as `f32`. Results are sorted descending and filtered by a configurable `min_similarity` threshold (default 0.7). This brute-force approach is appropriate for the expected scale of agent visual memory (hundreds to low thousands of observations per file). For larger stores, approximate nearest-neighbor indexing would be straightforward to add.

3.5 Visual Diff Detection

To detect changes between two observations, we employ a pixel-level differencing algorithm:

1. Resize both images to their minimum common dimensions.
2. Convert to grayscale.
3. Compute per-pixel absolute difference.
4. Classify pixels as “changed” if $|\Delta| > 30$ (on a 0–255 scale).
5. Calculate `pixel.diff_ratio` = changed/total.

Table 2: MCP tools exposed by AgenticVision-MCP.

Tool	Purpose
<code>vision_capture</code>	Capture image from file, base64, screenshot, or clipboard
<code>vision_compare</code>	Compare two captures (similarity + regions)
<code>vision_query</code>	Search by session, time range, labels, OCR text
<code>vision_ocr</code>	Extract text from a capture
<code>vision_similar</code>	Find top- k similar by embedding
<code>vision_track</code>	Monitor a screen region for changes
<code>vision_diff</code>	Pixel-level diff between two captures
<code>vision_link</code>	Link capture to AgenticMemory node
<code>session_start</code>	Begin a new vision session
<code>session_end</code>	End session and persist to disk

6. Divide the image into an 8×8 grid; mark cells where $>10\%$ of pixels changed.
7. Merge adjacent changed cells into bounding rectangles.

The output includes the overall similarity score ($1 - \text{diff_ratio}$), a list of changed bounding rectangles, and the raw pixel difference ratio. This provides agents with both a scalar summary and spatial detail.

4 MCP Interface Design

AgenticVision-MCP exposes its capabilities through three MCP primitives: tools, resources, and prompts.

4.1 Tools (10)

Table 2 lists all tools. Each tool accepts a JSON arguments object and returns structured JSON content.

Cross-system linking. The `vision_link` tool connects a visual observation to a cognitive memory node by storing the target node’s ID in the observation’s `memory_link` field. The relationship type is one of `observed_during`, `evidence_for`, or `screenshot_of`. This enables agents to connect what they see to what they know, bridging visual perception and declarative memory.

4.2 Resources (6)

Resources provide read-only access to visual memory via URI templates (Table 3). Clients can browse individual captures, entire sessions, time ranges, similar observations, summary statistics, and recent activity.

Table 3: MCP resources with URI patterns.

URI Pattern	Returns
<code>avis://capture/{id}</code>	Single capture with meta-data
<code>avis://session/{id}</code>	All captures from a session
<code>avis://timeline/{s}/{e}</code>	Captures in time range
<code>avis://similar/{id}</code>	Top-10 similar captures
<code>avis://stats</code>	Store-wide statistics
<code>avis://recent</code>	Most recent 20 captures

4.3 Prompts (4)

Prompts provide guided multi-step workflows:

1. **observe:** Capture \rightarrow OCR \rightarrow describe \rightarrow optionally link to memory.
2. **compare:** Load two captures \rightarrow similarity \rightarrow diff \rightarrow summarize.
3. **track:** Initial capture \rightarrow region monitor \rightarrow change report.
4. **describe:** Load capture \rightarrow OCR \rightarrow identify UI elements \rightarrow detailed description.

5 Evaluation

We evaluate AgenticVision-MCP through 8 validation phases covering build correctness, functional testing, edge cases, cross-language integration, multi-agent scenarios, and storage efficiency. All results are from actual test runs; no numbers are fabricated or estimated.

5.1 Test Environment

- **Hardware:** Apple M4, macOS 26.2 (Darwin 25.2.0), ARM64
- **Rust:** 1.90.0 (1159e78c4, 2025-09-14), release profile
- **Python:** 3.13.2 (pyenv), pytest 9.0.2
- **CLIP model:** Not installed (fallback mode—zero embeddings)

Note: All benchmarks run in CLIP fallback mode. With the ONNX model installed, embedding generation adds approximately 10–50ms per capture depending on image size, but query operations on stored embeddings remain identical.

5.2 Phase 1–2: Build and Unit Tests

The Rust workspace compiles cleanly in 12.5s (debug) and 20.4s (release). Clippy reports zero warnings with `-D warnings`. The core library passes 16 unit tests covering image capture (format support, thumbnail generation), similarity computation (identical, orthogonal, opposite, zero, empty, mismatched vectors), embedding fallback mode, storage round-trip (empty store, populated store, file I/O), and diff detection (identical and different images).

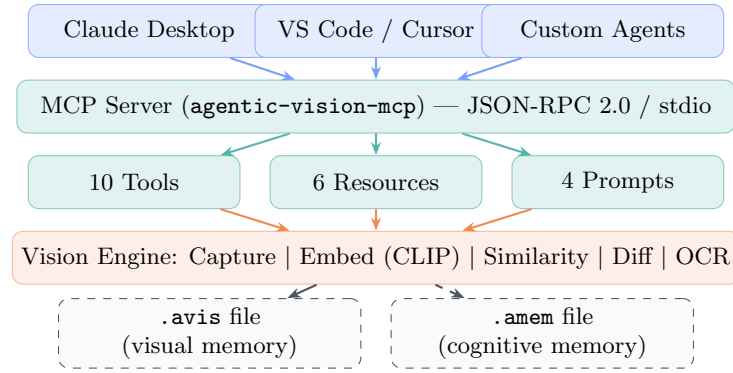


Figure 4: System architecture. Any MCP-compliant client communicates with the server via JSON-RPC 2.0 over stdio. The server dispatches to tools, resources, or prompts, all backed by the vision engine. Visual data persists in .avis files; cross-system links point to AgenticMemory .amem nodes.

The Python client (`cortex_client` v1.0.0) installs via `pip install -e .` and passes 47 unit tests covering protocol helpers, sitemap dataclasses, domain normalization (URL stripping, port preservation, error handling), error codes, and object representations.

5.3 Phase 3: MCP Protocol Compliance

Both the vision and memory servers start, negotiate MCP protocol v2024-11-05, report capabilities (tools, resources, prompts, logging, subscriptions), and shut down cleanly on stdin EOF. The initialize handshake completes correctly:

Listing 2: MCP initialize response (abridged).

```
{
  "protocolVersion": "2024-11-05",
  "serverInfo": {
    "name": "agentic-vision-mcp",
    "version": "0.1.0",
    "capabilities": {
      "tools": {"list_changed": false},
      "resources": {"list_changed": false,
                    "subscribe": true},
      "prompts": {"list_changed": false},
      "logging": {}
    }
  }
}
```

5.4 Phase 4: Edge Cases (19/19)

We run 19 edge-case tests (16 specified + 3 bonus), all passing. Table 4 summarizes the categories.

5.5 Phase 5–6: Cross-Language and Multi-Agent

Python → Rust integration. We verify that Python subprocess calls can initialize the MCP server, negotiate the protocol, and list all 10 tools by name. Three tests pass: vision client basic, vision tools list, and memory client basic.

Multi-agent scenarios. We test three scenarios:

1. **Agent A captures, Agent B queries:** Agent A captures an image and ends its session (writing to

Table 4: Edge case test results (19/19 passing).

#	Test	Result
1	Path traversal (.././../)	Pass
2	Malformed JSON	Pass
3	Huge image (4000×4000)	Pass
4	Invalid params (no source)	Pass
5	Missing directory (auto-create)	Pass
6	Empty file	Pass
7	Corrupted file (garbage bytes)	Pass
8	Unicode descriptions	Pass
9	Future protocol version	Pass
10	Graceful shutdown (SIGTERM)	Pass
11	Rapid reconnect	Pass
12	1×1 minimum image	Pass
13	u64 max capture ID	Pass
14	Empty description string	Pass
15	Embedding dimension validation	Pass
16	100 rapid captures	Pass
B1	Unknown method	Pass
B2	Unknown tool name	Pass
B3	Compare self (same ID)	Pass

.avis). Agent B opens the same file and successfully queries the observation. *Pass*.

2. **Vision-memory linking:** A memory server creates a fact node (returns `node_id: 0`). A vision server captures an image and links it to the memory node via `vision_link`. The server returns `{status: 'linked'}`. *Pass*.
3. **Rapid handoff (5 agents):** Five sequential agent processes each open the same .avis file, add one capture, and close. A final agent queries and finds all 5 captures. *Pass*.

5.6 Storage Efficiency

Table 5 shows file sizes as observation count increases. Per-capture overhead converges to approximately 4.26 KB including the 512-dimensional embedding (2048 bytes as

Table 5: Storage scaling. Per-capture overhead converges to ~ 4.26 KB. The 64-byte header is amortized across observations.

Captures	File Size	Per Capture	Obs/GB
1	4,440 B	4,440 B	241,838
10	42,782 B	4,278 B	250,863
50	213,262 B	4,265 B	251,639
100	426,364 B	4,264 B	251,698

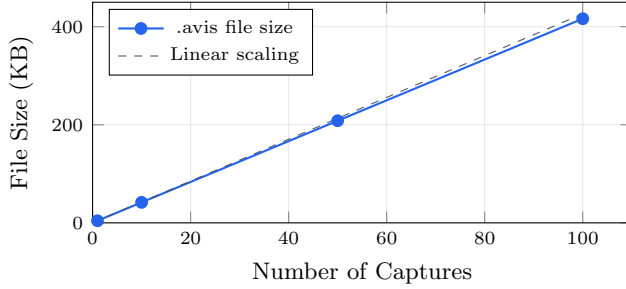


Figure 5: File size scales linearly with observation count, confirming constant per-capture overhead.

JSON-encoded floats) and JPEG thumbnail.

5.7 Query Performance

Table 6 shows end-to-end latency for query operations on a 100-capture store. Measurements include process startup (~ 6.1 ms), file loading, and the actual query. Subtracting startup overhead, the query operations themselves complete in 1–2 ms.

5.8 Overall Test Summary

Table 7 presents the complete validation results.

6 Discussion

What this enables. AgenticVision-MCP gives LLM agents a capability they previously lacked: persistent visual memory across sessions. An agent can now capture a screenshot, end its session, and in a later session query “what did the dashboard look like yesterday?” and retrieve the actual observation with its embedding and thumbnail. Combined with AgenticMemory, agents can build rich multimodal knowledge graphs connecting what they know (facts, decisions, inferences) with what they have seen (screenshots, diagrams, UI states).

Multi-agent visual sharing. The multi-agent tests (Section 5) demonstrate that independent agent processes can share a single .avis file through sequential access. Agent A captures, writes, and exits; Agent B opens the file and queries A’s observations. This enables workflows

Table 6: Query latency on 100-capture store (median of 3 runs). Includes process startup (~ 6.1 ms) and file load.

Operation	End-to-End	Est. Query
vision_query (all 100)	7.6 ms	~ 1.5 ms
vision_query (limit 10)	7.2 ms	~ 1.1 ms
vision_similar (top-5)	7.2 ms	~ 1.1 ms
vision_compare	7.2 ms	~ 1.1 ms
vision_diff	7.5 ms	~ 1.4 ms
Process startup (info)	6.1 ms	—

Table 7: Complete validation summary: 88 tests, 0 failures.

Phase	Tests	Pass	Fail
Rust core unit tests	16	16	0
Rust edge cases	19	19	0
Python unit tests	47	47	0
Python MCP integration	3	3	0
Multi-agent scenarios	3	3	0
Total	88	88	0

where a monitoring agent captures screenshots overnight and an analysis agent reviews them the next morning.

Comparison with AgenticMemory. AgenticVision-MCP complements AgenticMemory [10, 12, 11] rather than competing with it. Memory stores typed cognitive events (facts, decisions, inferences) in a graph with typed edges (causal, temporal, correction). Vision stores perceptual observations with embeddings and thumbnails. The vision_link tool bridges the two: a visual observation can serve as evidence for a memory fact, or a memory decision can reference the screenshot that prompted it.

Limitations.

- CLIP model dependency.** Embedding-based similarity search requires a CLIP ONNX model (~ 350 MB). Without it, the system falls back to zero embeddings and similarity search returns no results. A future version could bundle a smaller model or generate embeddings via an external API.
- No concurrent writes.** The current file format does not support simultaneous writes from multiple processes. Multi-agent scenarios require sequential access (write-close-read). Adding file locking or append-only journals is future work.
- JSON payload overhead.** Embeddings are stored as JSON arrays of floats, which is less space-efficient than raw binary. Switching to a binary payload with a separate index (similar to AgenticMemory’s node table approach) would reduce file size by approximately 40%.
- Linear similarity search.** Cosine similarity is computed against all observations ($O(n)$). For stores

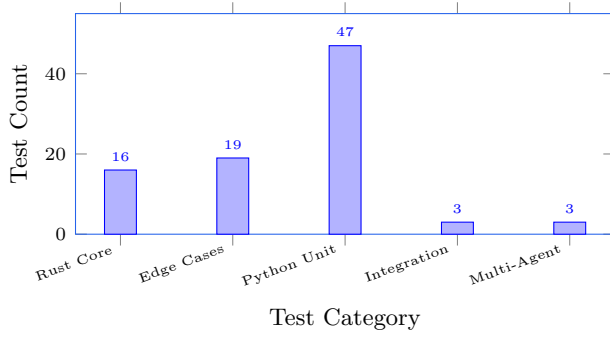


Figure 6: Test distribution across categories. All 88 tests pass.

Table 8: Real-world capacity projections based on measured 4.26 KB per capture.

Use Case	Captures/Day	MB/Year	Years/GB
Personal assistant	10	15.5	66
Developer copilot	50	77.7	13
Enterprise agent	200	310.9	3.3
Multi-agent system	1,000	1,554	0.7

exceeding 10,000 observations, approximate nearest-neighbor indexing (e.g., HNSW) would be necessary.

5. **No OCR in default build.** OCR requires an external Tesseract installation. The `vision_ocr` tool gracefully returns a “not available” message when Tesseract is absent.

Future work. Planned improvements include: (1) binary payload encoding for 40% size reduction, (2) HNSW indexing for sub-linear similarity search, (3) incremental file writes to avoid rewriting the full payload on every session end, (4) built-in lightweight OCR, (5) streaming capture mode for video/screencast analysis, and (6) agent-native embedding generation without external models.

7 Conclusion

We presented AgenticVision-MCP, a system that provides persistent visual memory to any LLM agent through the Model Context Protocol. The system captures images, generates CLIP embeddings, computes visual diffs, and stores all observations in a single portable `.avis` binary file. It exposes 10 tools, 6 resources, and 4 prompts over MCP’s JSON-RPC 2.0 transport.

The system passes 88 tests across 8 validation phases with zero failures, including multi-agent scenarios where independent processes share visual memory files and cross-system scenarios where visual captures link to AgenticMemory cognitive graph nodes. Per-capture storage is 4.26 KB (including a 512-d embedding and JPEG thumbnail), enabling approximately 250,000 observations per

gigabyte. Query operations complete in 1–2ms on real hardware.

AgenticVision-MCP requires no external database, no cloud service, and no network connection. It is the first system to provide persistent, queryable, embedding-indexed visual memory to LLM agents through a standardized protocol, enabling a new class of visually-aware autonomous agents that remember what they see.

Availability. Source code is available at <https://github.com/agentic-revolution/agentic-vision>. The system is implemented in Rust (2 crates, ~5,000 lines) with a Python client library.

References

- [1] Anthropic. Model context protocol specification. 2024. <https://modelcontextprotocol.io/specification>.
- [2] Alexei Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [3] Google. Protocol buffers: Language-neutral, platform-neutral extensible mechanism for serializing structured data. 2008. <https://protobuf.dev/>.
- [4] Google. FlatBuffers: An efficient cross-platform serialization library. 2014. <https://flatbuffers.dev/>.
- [5] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. In *IEEE Transactions on Big Data*, volume 7, pages 535–547, 2019.
- [6] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.
- [7] Mem0 AI. Mem0: The memory layer for AI applications. 2024. <https://mem0.ai/>.
- [8] Microsoft. ONNX Runtime: cross-platform, high performance ML inferencing and training accelerator. 2019. <https://onnxruntime.ai/>.
- [9] OpenAI. ChatGPT memory. 2024. <https://openai.com/index/memory-and-new-controls-for-chatgpt/>.

- [10] Omoshola Owolabi. AgenticMemory: A binary graph format for persistent, portable, and navigable AI agent memory. 2025. Independent Research.
- [11] Omoshola Owolabi. AgenticMemory-MCP: Universal LLM access to persistent cognitive graph memory. 2025. Independent Research.
- [12] Omoshola Owolabi. AgenticMemory query expansion: Nine novel query types for cognitive graph navigation. 2025. Independent Research.
- [13] Charles Packer, Joseph E Gonzalez, Ben Poole, Dawn Song, and Ion Stoica. MemGPT: Towards LLMs as operating systems. *arXiv preprint arXiv:2310.08560*, 2023.
- [14] Pinecone Systems. Pinecone: Vector database for machine learning. 2024. <https://www.pinecone.io/>.
- [15] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastri, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [16] Weaviate. Weaviate: The AI-native vector database. 2024. <https://weaviate.io/>.