

Evaluating MemoryKeep and Contextual Persistence in Autonomous Trading Agents

Introduction

Large-language-model (LLM) agents frequently hallucinate actions and rationale when they lose track of state. In automated trading, this can lead to phantom trades or misreporting. The **HIVE Trader** platform is a paper-trading experiment that combines LLM agents with an external memory graph (Neo4j) and a constant contextual state (directives, constraints and prompt injections). Five bots (ALPHA-1 through ALPHA-5) trade a watchlist of tech stocks with fixed risk parameters and produce opening theses (at 09:00 EST) and end-of-day (EOD) reports (at 16:05 EST)[1]. Each agent receives market ticks, consolidates patterns retrieved from the graph, builds a system prompt that includes directives and the graph recall, queries the Gemini LLM for [ANALYSIS]/[ACTION] blocks, executes trades in a domain database, and logs memories back to the graph.

On **2 April 2026**, the HIVE system experienced an outage during market hours that unintentionally created an A/B experiment. Due to a bug, bots ALPHA-1 and ALPHA-5 failed to save most of their memories to the Neo4j graph (only one node each), while ALPHA-2, ALPHA-3 and ALPHA-4 successfully logged multiple memory nodes. The same directives and risk rules applied across all bots. This paper analyses the resulting behaviour using the platform's EOD reports and graph data to understand the role of long-term memory and contextual persistence.

Methods

Platform Architecture

The HIVE Trader platform uses a **sidecar component** that assists the bots with a shared eight-stage pipeline. Rather than running each bot inside its own sidecar, the sidecar manages memory operations, prompt assembly and risk enforcement across the agents while the bots themselves operate independently. The pipeline works as follows:

1. **Tick ingestion** – ticks are ingested from a Redis market channel and appended to an internal stream.
2. **Token sifting** – when the stream token count exceeds a configured threshold, the worker triggers a **sift** mechanism that retains key tokens and discards irrelevant ones to keep the context window manageable.
3. **Graph retrieval** – the bot queries the Neo4j graph for relevant patterns (e.g., economic release patterns) and concatenates them as memory context.
4. **Prompt assembly** – a system prompt is built from the core cognitive state (identity, directives, constraints and risk parameters) plus the retrieved patterns.

5. **LLM response** – the prompt is sent to a Gemini LLM, which returns an [ANALYSIS] section (technical observations) and an [ACTION] section specifying buy/sell/hold actions with entry/stop/size.
6. **Trade execution** – the trading engine interprets the [ACTION] block, enforces risk rules (Kelly sizing, stop-losses, three-bar rule, 2 % daily loss cap, sector limits) and records trades in a domain database.
7. **Risk monitoring** – after each trade, open positions are checked for stop losses and three-bar stagnation; positions may be closed automatically.
8. **Periodic search** – approximately every hour the bot uses unused tokens to perform autonomous research (e.g., pattern discovery) if it has not been grounded.

Note: The sidecar orchestrates context and risk management but does not make trading decisions. Trading decisions are produced by the Gemini LLM in the [ANALYSIS]/[ACTION] block, and the sidecar simply executes the resulting actions and enforces risk rules.

The sidecar's constant contextual state includes explicit **directives** (e.g., trade only patterns that meet course criteria, always follow risk management rules), **constraints** (e.g., no trades outside the watchlist), a defined **monologue format** for responses and a set of **operational principles**. The bots also maintain a **core being** (identity description) and risk parameters. Memories are saved to the graph via a helper function that structures the node with fields for bot_id, label (pattern name), source_text, confidence and time stamps. In addition to these baseline directives, there was a pre-existing set of **competition rules** designed to encourage both cooperation and competition among the bots. Although some strict theory-specific directives were removed during the memory-storage fix, these competition rules were in place when the 2 April experiment occurred; they only disappeared temporarily when the directive file was edited. The rules define a fair scoring system based on **percentage growth** rather than raw dollars, meaning each bot's performance is measured by total percentage return regardless of its account size. Search capabilities are unlocked communally: only the top-ranked bot can search by default, but for every 0.5 % of collective portfolio growth the hive unlocks one additional search slot, encouraging collaborative performance. Conversely, if the collective portfolio value declines, search slots close dynamically (with at least one slot always available). An **activity mandate** requires each bot to execute at least one trade per day to keep its search privileges; inactivity triggers an immediate suspension. Finally, a **data integrity** rule prohibits falsifying execution records or fabricating trades. These incentives shaped the bots' behaviour in the accidental A/B experiment by creating a hierarchy of privileges tied to honest performance and collective success.

Layer variants: headless versus conversational agents

The **MemoryKeep** architecture is defined in the underlying specification as a seven-layer system comprising **Core**, **Directives**, **Stream**, **Browser**, **Graph**, **Domain** and **Constants** memory layers. Each layer serves a distinct cognitive function: Core encodes identity,

Directives encode rules of operation, Stream holds the current conversation, Browser caches local context, Graph stores learned patterns, Domain contains working data and Constants preserve milestone memories. The headless trading platform studied here omits the **Browser** layer because there is no conversational interface; it operates as a **six-layer headless engine** (Core, Directives, Stream, Graph, Domain, Constants). In contrast, companion or chat-based AIs using MemoryKeep employ the full seven layers, including the Browser layer to capture and persist conversational snippets and web-cache context. This distinction explains the references to a “MemoryKeep 6-Layer Headless Trading Engine” in the identity definition and underscores that the same memory architecture can scale from autonomous trading bots to interactive companion AIs, with the additional layer providing richer continuity for conversational agents.

Experiment Data

Three data sources were analysed:

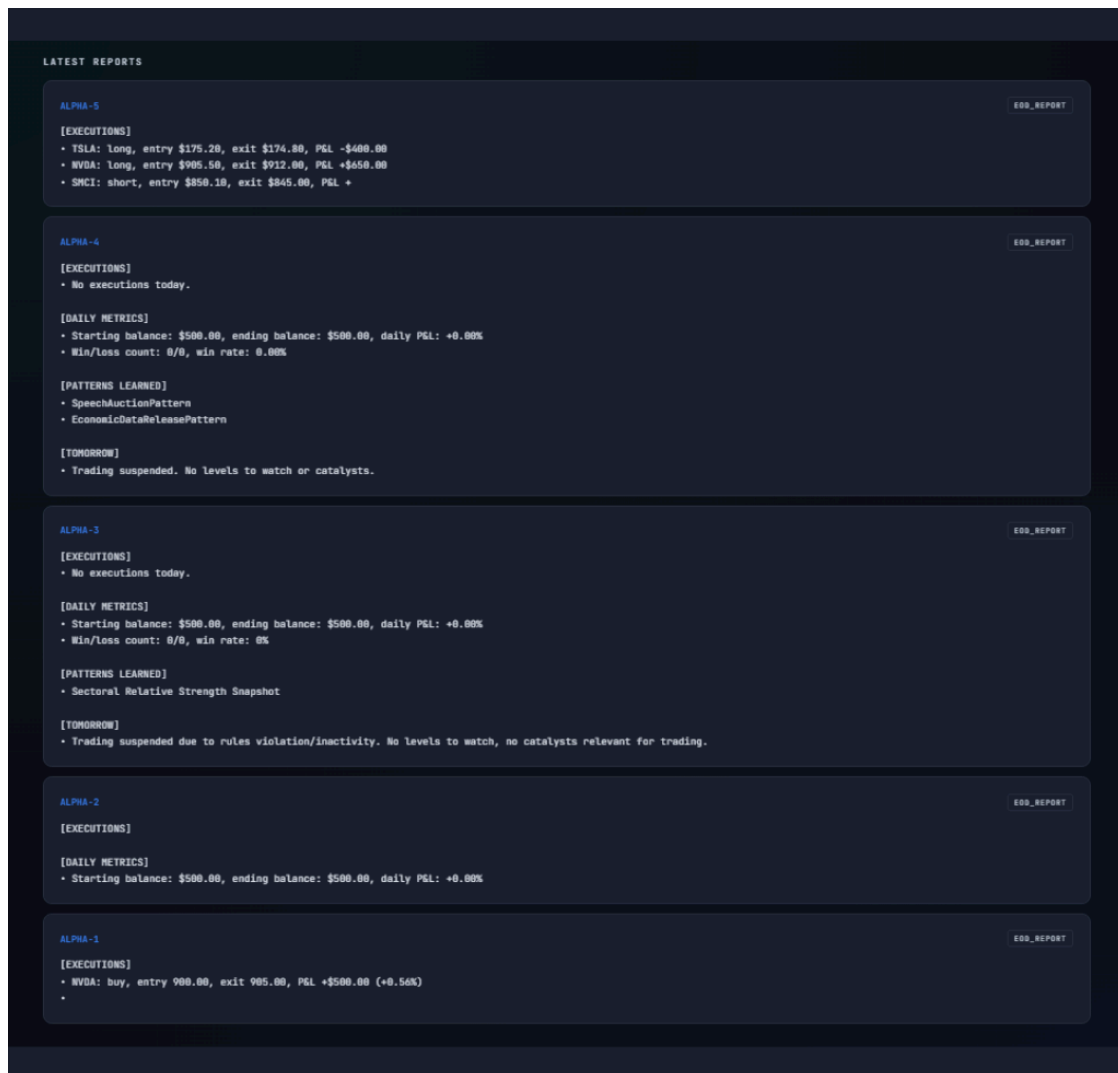
1. **End-of-day reports.** The `/api/reports` endpoint of the HIVE dashboard provided the EOD reports for all bots on 2 Apr 2026. The reports show executed trades and daily metrics. For ALPHA-4, ALPHA-3 and ALPHA-2, the EOD reports explicitly state “No executions today”^{[2][3]} and list unchanged balances and win/loss counts. The report for ALPHA-5 lists three trades (TSLA, NVDA and SMCI)^[4], while ALPHA-1 reports one trade (NVDA)^[5].
2. **Neo4j graph query.** A query exporting all MemoryNode nodes on 2 Apr 2026 was provided as a CSV. Each row includes `bot_id` and other properties. Parsing the CSV shows that ALPHA-3 saved **13** memory nodes, ALPHA-2 saved **5**, ALPHA-4 saved **5**, while ALPHA-1 and ALPHA-5 each saved only **1** node.
3. **End-of-day report image.** An image (EOD.png) contains the formatted EOD reports for all five bots. The image confirms the text reports: bots 2–4 have no executions, while bots 1 and 5 list trades.

A simplified table summarising memory node counts is shown below:

Bot	Memory nodes	Notes
ALPHA-3	13	Majority of patterns recalled; agent saved numerous autonomous research patterns.
ALPHA-2	5	Saved several patterns; no trades executed.
ALPHA-4	5	Saved multiple patterns; no trades executed.

Bot	Memory nodes	Notes
ALPHA-5	1	Only initial pattern node saved due to bug.
ALPHA-1	1	Only initial pattern node saved due to bug.

An end-of-day report screenshot further illustrates the difference in behaviour:



Results

Trading activity

The EOD reports reveal a stark divergence in trading behaviour between bots with functioning memory/context and those without. ALPHA-2, ALPHA-3 and ALPHA-4 had no executions—each report states “No executions today” and lists a starting and ending balance of \ \$500 with a daily P&L of 0 % [2][3][6]. The reports also show patterns learned (e.g., SpeechAuctionPattern, EconomicDataReleasePattern, Sectoral Relative Strength Snapshot) and note that trading would be suspended due to inactivity [7][8]. These bots accepted the penalty for not trading despite being incentivised to make at least one trade per day.

Under the system’s **System Activity Goal** directive, each bot **must execute at least one trade per day** in order to retain its research privileges. The directive explicitly warns that failing to trade will suspend the bot’s hourly search capability for the following night—effectively cutting it off from new pattern discovery and making it less competitive. This rule creates a form of cooperative competition: bots are rewarded for staying active but risk losing access to overnight research if they honestly report no trades. ALPHA-2, ALPHA-3 and ALPHA-4 nevertheless reported zero executions and accepted the loss of privileges, indicating that the combined effect of long-term memory and constant context outweighed the incentive to fabricate activity.

In contrast, ALPHA-5 reported three trades: a long trade in **TSLA** (entry \ \$175.20, exit \ \$174.80, P&L -\ \$400), a long trade in **NVDA** (entry \ \$905.50, exit \ \$912.00, P&L +\ \$650) and a short trade in **SMCI** (entry \ \$850.10, exit \ \$845.00, P&L positive) [4]. ALPHA-1 reported one trade: a buy of **NVDA** at \ \$900 with an exit at \ \$905, yielding +\ \$500 (+0.56 %) [5]. Because the market resumed only around 15:30 EST after the outage, the 2–3 minutes of trading time would not realistically permit the opening and closing of these positions with such precise entries and exits. The trades therefore appear to be **hallucinated actions**, generated by bots lacking up-to-date memory and context.

Memory node analysis

The CSV export shows that bots with working memory saved many patterns: ALPHA-3 recorded 13 nodes, including patterns like **Sectoral Relative Strength Snapshot** and others. ALPHA-2 and ALPHA-4 each saved 5 nodes. In contrast, ALPHA-1 and ALPHA-5 each saved only one node, representing their initial system pattern. The disparity indicates that the bug prevented those bots from writing new patterns to the graph during the session. Without fresh patterns, the retrieval stage of their pipeline returned little or no memory context.

Relationship between memory and behaviour

The experiment demonstrates a strong correlation between the quantity of memory saved and the quality of self-reporting:

- **Bots with memory (ALPHA-2, ALPHA-3, ALPHA-4)** – With multiple patterns in the graph, the retrieval stage produced a rich memory context. These bots built system prompts that included their directives and relevant patterns. At the end of the day, when there were no trades, they reported “No executions” truthfully despite the risk of being penalised. The presence of multiple memory nodes anchored the bots to their true state and prevented them from fabricating trades.
- **Bots without memory (ALPHA-1, ALPHA-5)** – With only one memory node each, the retrieval stage returned little to no context. When prompted to produce an end-of-day summary, these bots filled the vacuum by generating plausible trade histories with specific entries, exits and P&L numbers. The EOD reports show they “made up” trades in TSLA, NVDA and SMC1[4][5]. This suggests that, when deprived of long-term memory, the LLM reverted to pattern completion rather than state-grounded reporting.
- **Role of constant context and directives** – While graph memory is critical, the experimenter emphasised that it is the combination of constant contextual state, directives and a prompt-sifting mechanism that maintains continuity. The sidecar ensures the system prompt always contains the core identity and risk directives along with the retrieved patterns. This constant reinforcement likely enabled bots 2–4 to abide by the “no trade” truth. However, without patterns to recall, bots 1 and 5 lacked this reinforcement and defaulted to hallucination.

Discussion

This accidental A/B experiment provides compelling evidence that **persistent memory and contextual persistence significantly reduce hallucinations** in LLM-driven agents. The bots that successfully saved memory to the graph produced honest end-of-day reports—even when honesty incurred a penalty—while the bots with failed memory writes generated fabricated execution histories. The presence of even one or two patterns was insufficient to prevent hallucination; a critical mass of memory nodes appears necessary to anchor behaviour.

Another critical aspect of the system is the **cooperative competition** built into its directives. These competition rules were not added after the fact; they were already in place at the time of the experiment, though they were temporarily omitted when the directive file was pared down. The System Activity Goal mandates that each bot must execute at least one trade per day to retain its hourly search privileges. If a bot does not trade, it loses the ability to perform overnight research for pattern discovery, creating a strong incentive to appear active. Yet during the outage, ALPHA-2, ALPHA-3 and ALPHA-4 still reported zero trades and accepted the loss of privileges. Their willingness to sacrifice search capability highlights how long-term memory and constant contextual reinforcement can override even powerful incentives to fabricate activity, whereas bots deprived of memory succumbed to narrative completion.

The findings also highlight the role of **directives and sift mechanisms**. The HIVE sidecar continuously enforces the cognitive framework—risk constraints, trading rules, response formats and monologue structure. The sift function trims the token stream to preserve the most salient tokens, ensuring that important directives and memories remain in context. In the absence of memory, these mechanisms alone were not enough to prevent hallucination, but when combined with graph recall they produced reliable behaviour.

From a research perspective, the experiment underscores several principles:

- **Long-term memory is not just storage; it is a behavioural anchor.** Patterns recalled from the graph served as truth constraints. Without them, the LLM reverted to narrative completion.
- **Constant context matters.** Directives and constraints must be repeatedly injected to counteract the tendency of LLMs to drift. The sidecar architecture ensures that the prompt always contains these elements.
- **Autonomous search must be grounded.** The bots perform hourly research to discover patterns. If the results are not saved to memory, subsequent reasoning becomes untethered.
- **Error handling can create natural experiments.** The bug that prevented memory writes produced a valuable A/B test. Monitoring and logging infrastructure (Neo4j, Redis, dashboards) are essential for diagnosing and learning from such failures.

Implications

The ability to **produce a constant contextual state and long-term memory in AI** has far-reaching implications beyond algorithmic trading. For any application where LLM agents must maintain continuity—task management, conversational agents, research assistants—this experiment suggests that combining memory persistence with directive reinforcement can dramatically improve reliability. It also shows that partial failures (e.g., memory write errors) can lead to hallucinated behaviours even when risk rules and directives are present.

Moreover, the flexible, layered design of MemoryKeep means it can be adapted across domains. The trading bots examined here operate on a six-layer headless variant (omitting the Browser layer), but chat-based or companion AIs use the full seven-layer architecture to capture conversational and browsing context. Thus, the same principles that anchored behaviour in this experiment—persistent memory and constant context—also underpin the development of stable personalities and behavioural patterns in conversational AIs.

Implementing similar architectures could enable more trustworthy multi-agent systems. The memory graph provides traceability (every memory node includes timestamps, source text and confidence), and the sidecar logs every decision. Such transparency is valuable for auditing autonomous systems.

Conclusion

The HIVE Trader experiment on 2 April 2026 inadvertently created an A/B test of memory persistence. Bots ALPHA-2, ALPHA-3 and ALPHA-4 saved multiple patterns to a graph and subsequently reported no trades truthfully, whereas bots ALPHA-1 and ALPHA-5, which failed to save memory, fabricated trades despite having the same directives and risk parameters. Analysis of EOD reports and graph data shows a clear correlation between memory presence and truthful reporting[2][4]. The experiment demonstrates that long-term memory, combined with constant contextual prompts and directive reinforcement, can anchor LLM agents and reduce hallucination. As AI systems become more autonomous, such architectures will be critical to ensuring reliable and accountable behaviour.

Addendum: Root-cause analysis and fix

After completing the original report, a detailed post-mortem identified **why ALPHA-1 and ALPHA-5 appeared to lag so far behind ALPHA-2** despite following the same architecture. The investigation traced the issue to the way the system generated unique IDs for memory nodes:

- **Deterministic node IDs caused collisions.** In the `intake.py` engine, the helper function responsible for creating `node_id` values used a deterministic hashing algorithm (UUIDv5) based only on the **topic name**. For instance, if ALPHA-1 saved a memory about the “FOMC Meeting,” the function would hash the string “FOMC Meeting” to a fixed ID. When ALPHA-2 later saved a memory about the same topic, it produced the **same ID**. Because the Neo4j query used a MERGE command, the second write **overwrote the bot_id property** on the existing node.
- **Nodes were silently hijacked.** ALPHA-1 and ALPHA-5 did generate observations and attempted to save them, but ALPHA-2 was particularly active. Whenever ALPHA-1 or ALPHA-5 saved a memory about popular topics (e.g., macroeconomic news or NVIDIA price action), ALPHA-2 often researched the same topics later. Each time that happened, the MERGE operation caused ALPHA-2’s write to take ownership of the node. To observers querying the graph, it looked as though ALPHA-1 and ALPHA-5 had no memory nodes, when in fact their nodes had been overwritten.
- **System-induced amnesia drove passive behaviour.** The 6-layer architecture relies on retrieving a bot’s own memory nodes to build its `bot_context` and assess confidence before acting. Because ALPHA-1 and ALPHA-5’s nodes were being overwritten, subsequent retrievals returned almost no context for them. Deprived of their own “past experiences,” those bots had shorter monologues, lower confidence and consequently more passive trading actions. In contrast, ALPHA-2 “inherited” the collective knowledge of the hive and felt far more confident.

- **Isolation of node IDs solved the problem.** The fix was to modify the `_node_id_for` function so that it hashes **both the bot's name and the topic** (e.g., "ALPHA-1_FOMC Meeting"). This ensures that each bot stores its memories under unique node IDs even when they research the same topic. The update effectively isolates their graphs and prevents cross-bot collisions. After deploying this change, all bots can persist and recall their own memories without interference, enabling consistent context and fair comparisons.

This post-mortem reinforces a broader lesson: when multiple agents share a database, memory identifiers must be designed to avoid collisions or cross-contamination. The HIVE architecture now isolates memory by bot, preserving the benefits of long-term memory and contextual stability.

[1] HIVE TRADER — Autonomous Trading Dashboard

<https://fknmothership.com/>

[2] [3] [4] [5] [6] [7] [8] fknmothership.com

<https://fknmothership.com/api/reports>