

How Two Databases Quietly Break Your RBAC and GDPR Story

Second Order Labs™ · July 2026

ABSTRACT

A vector workload that dropped to \$200/month on pgvector exposed the real cost of two-database architectures: the operational and compliance drag that never shows up on an invoice.

Keywords: Vector Databases, pgvector, RAG, Hybrid Search, Data Architecture, AI Infrastructure

“The standalone vector database isn't dead. But it's being cornered into the billion-scale niche where most teams will never operate.”

— Author of Vector Databases Are Dying, Vector Databases Are Dying. Here's the Production Evidence.

I. The \$200 Migration That Exposed a \$50M Assumption

A mid-sized fintech engineering team processing 50 million daily transactions ran the same retrieval workload on a standalone vector database and on pgvector. Same queries. Same results. The standalone bill collapsed to roughly \$200 a month after migration.^[1] That number tells the story everyone already knows: standalone vector databases are expensive at scale. But it buries a more dangerous reality. The migration was inevitable because the two-database architecture was quietly breaking things that never showed up on an invoice.

The cost story is well documented. Vector databases can consume 40 to 50 percent of total application cost in production, second only to LLM API calls.^[2]

Metadata filtering silently multiplies read-unit consumption, turning a one-unit semantic query into five or ten.^[1] Per-query billing works at prototype scale and fails at production scale.^[3] What no one prices is the compliance blast radius of running relational data and vector data as two independent stateful systems.

Call it **Dual-State Drag**: the compounding operational tax, which also drags down security and velocity, paid by teams keeping source-of-truth records and their vector representations in two databases.

These systems must be kept consistent but never actually can be. Dual-State Drag is the hidden second-order cost that margin-cannibalization coverage misses. It drives the shift toward converged architectures.

II. Why Two Databases Can Never Stay Consistent

Dual-state architectures fail because consistency between a relational store and a vector store is asynchronous by construction. Every write to the primary database triggers a separate, delayed operation to regenerate and reinsert an embedding. That gap is where correctness dies.

When a new product lands in your primary database, you face what one Google Cloud engineer described as a "frantic, asynchronous scramble to create its vector fingerprint and stuff it into the vector database."^[4]

That scramble is not an edge case. It runs on every database mutation forever. Teams start with a stand-alone vector database for retrieval. Then they discover that embeddings and permissions must stay perfectly synced with operational records.^[5] The moment they discover it, they own a distributed-consistency problem they never chose to build.

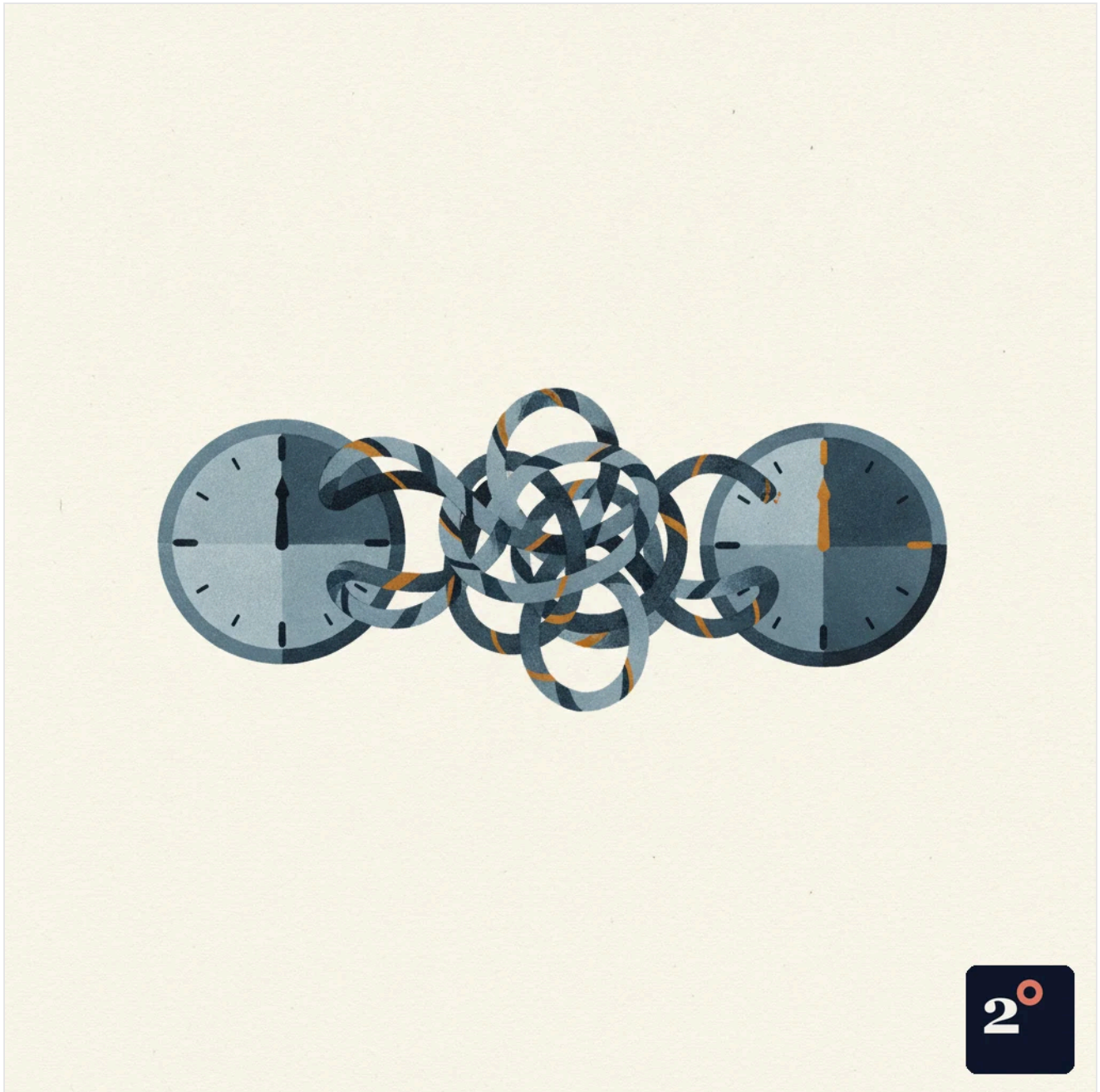


Figure 1: Dual-state architectures run on two clocks that never quite agree.

The drag compounds beyond runtime. Coordinating schema migrations across two independent stateful databases turns routine CI/CD into a choreography problem. Running a separate database alongside operational data requires extra infrastructure and forces additional coordination when AI workloads need application context.^[6] Costs compound as the dataset scales. Every coordination point is a place where a deploy can leave the vector store describing a version of reality that no longer exists.

III. The Zombie Vector Trap Is a Compliance Problem, Not a RAM Problem

Stale embeddings left behind after operational data changes are a live compliance liability. When the primary record is updated or deleted but its embedding survives, you get a zombie vector. In an in-memory database, that zombie keeps occupying expensive RAM indefinitely.^[7] The RAM waste is the part practitioners built CLI tools to measure.^[7] The compliance exposure is the part no dashboard shows.

A zombie vector is a searchable copy of data that your system of record believes is gone. If a user exercises deletion rights and the relational row disappears, the embedding derived from that row can persist in the vector store. It remains fully retrievable, ready to be surfaced by an AI agent as though the deletion never happened. The asynchronous scramble that creates zombies on every update is the exact mechanism leaving deleted data alive in a second system. The insidious part is silence. A zombie vector produces no error. Instead, it returns a confident answer built on a deleted record, accumulating undetected as the dataset grows.

A zombie vector is a searchable copy of data your system of record swears it already deleted.

IV. RBAC and Retrieval Quality: Where Standalone Databases Break Enterprises

Role-based access control is the clearest example of Dual-State Drag. Permissions are nearly impossible to keep coherent across two separate databases. Your operational database enforces who can see what. Your vector database, sitting beside it, has its own idea. The two drift.

Permissions live as relationships and rules in the relational system. Embeddings in the vector store carry only whatever metadata you remembered to copy across, subject to the same asynchronous lag as everything else. Update a user's access in the primary system and there is a window where the vector store still answers as if the old permissions held. Enterprises are consolidating to fix this. Running pgvector inside enterprise PostgreSQL eliminates the separate vector database, killing the duplicated management overhead while keeping access control in one enforcement layer.^[8]

Figure 2: In observed production deployments, the vector database runs 40-50% of total application cost, second only to LLM API calls.^[2]

A retrieval-quality trap sits layered on top. Pure semantic search fails silently on the exact queries enterprises care most about. For precise, high-stakes lookups such as cancellation policies, SKUs, policy IDs, and dates, embedding models weight conceptual meaning over literal matches and miss the answer entirely.^[9] Hybrid search combining BM25 lexical matching with semantic embeddings is a production requirement for handling proper nouns and exact matches.^[11] A standalone vector database optimized purely for similarity fights you on the queries where being wrong is most expensive.

V. The Team You Split When You Split the Database

Dual-State Drag has a human cost the invoices never capture. Two databases bifurcate the engineering org that maintains them. Someone owns the relational schema and its migrations. Someone else owns the vector index and the pipeline required to scale it. The seam between them becomes a permanent coordination tax on feature velocity. Every feature that touches retrieval now requires both owners in the room. Every incident requires deciding which database lied.

The operational overhead of managing separate licensing and infrastructure negates the specialized performance benefits for teams operating under a billion vectors.^[6] That is not a throwaway line. The raw

retrieval speed you bought the standalone database for gets eaten by the cost of keeping two systems in step. The escape hatch is often simpler than the category admits. For some AI agent architectures, plain Markdown files as the primary memory system beat

managed vector databases on unit economics and latency.^[12] When the cheapest correct answer is a text file, the burden of proof shifts onto the specialized store.

Figure 3: Vector search is being absorbed as a feature; the standalone product retreats into a shrinking niche.

VI. The Billion-Scale Niche Is Real, and It Is Small

Standalone vector databases survive where the scale genuinely demands them. That territory is far narrower than three years of adoption implied. The honest boundary is billion-scale and beyond, where general-purpose databases still struggle and specialized indexing earns its keep. As datasets grow from millions to billions of records, purely in-memory indexes become financially unsustainable, forcing on-disk techniques like DiskANN.^[14] Pushing to 10 billion vectors requires distributed indexing and machinery the vast majority of enterprise use cases will never touch.^[15]

The incumbents are annexing even this frontier. Amazon OpenSearch now builds billion-scale vector databases in under an hour, indexing up to 10 times faster at a quarter of the cost using GPU acceleration.^[16] AWS S3 Vectors targets up to 90 percent cost re-

duction versus standalone solutions.^[17] Vector search has become a checkbox feature in cloud data platforms, not a standalone moat.^[18]

"The standalone vector database isn't dead. But it's being cornered into the billion-scale niche where most teams will never operate." Source: Vector Databases Are Dying. Here's the Production Evidence.

VII. What to Do Before Your Next Retrieval Decision

Start from converged and force the standalone case to justify itself. Default to keeping vectors next to the data in pgvector or an equivalent. Pgvector scale already delivers roughly 75 percent lower cost than Pinecone on comparable workloads while closing the performance gap.^[19] Reserve a standalone vector database for the genuine billion-scale regime, and be ruthless about whether you actually live there.

Before committing, run three checks:

Check	What to look for	Why it matters
Zombie vector audit	Embeddings of deleted records still retrievable	GDPR/deletion compliance, not just RAM
Access control layer	One enforcement point vs. two drifting copies	RBAC coherence across updates
Hybrid retrieval	BM25 + semantic, not semantic alone	Exact matches on SKUs and policy IDs fail silently without it

The larger shift is a reversal of the default. For three years the question was "which vector database?" The better question is whether the vector belongs in a separate database at all. Dual-State Drag is the reason

the answer is no. The standalone vector database is rapidly becoming a solution in search of a problem, waiting for a billion-scale dataset that most companies will never build.

KEY FINDINGS

Vector databases can eat 40-50% of total production application cost, second only to LLM API calls.

A single query with metadata filtering can consume 5 to 10 read units instead of one, and that cost stays hidden during sales calls.

Zombie vectors keep deleted records searchable in a second store, so a data-hygiene chore becomes a compliance exposure.

pgvectorscale runs comparable workloads at roughly 75% less cost than Pinecone while closing the performance gap.

Amazon OpenSearch builds billion-scale vector databases in under an hour with GPU acceleration, which commoditizes the last niche.

REFERENCES

- [1] Vector Databases Are Dying. Here's the Production Evidence. - Medium <https://medium.com/data-science-collective/vector-databases-are-dying-heres-the-production-evidence-8c17b54687e2>
- [2] How to Cut Vector Database Costs by Up to 80%: A Practical Milvus Optimization Guide <https://milvus.io/blog/how-to-cut-vector-database-costs-by-up-to-80-a-practical-milvus-optimization-guide.md>
- [3] Vector Database Pricing Comparison 2026: Real Cost Breakdown - RankSquire <https://ranksquire.com/2026/03/04/vector-database-pricing-comparison-2026/>
- [4] Google Cloud Database Digest: AlloyDB's ScaNN Vector Index Unifies Your Data & AI - Medium <https://medium.com/google-cloud/google-cloud-database-digest-alloydbs-scann-vector-index-unifies-your-data-ai-aug-22th-2025-3dc2eda8a345>
- [5] How to Choose a Database for an AI-Powered Product - CockroachDB <https://www.cockroachlabs.com/blog/database-for-ai-applications/>
- [6] Best Open Source Vector Databases 2026 & Comparison - Redis <https://redis.io/blog/best-open-source-vector-databases-comparison/>

- [7] I built a CLI to find "Zombie Vectors" in Pinecone/Weaviate - Reddit
- [8] pgvector: The Critical PostgreSQL Component for Your Enterprise AI Strategy - Percona <https://www.percona.com/blog/pgvector-the-critical-postgresql-component-for-your-enterprise-ai-strategy/>
- [9] Hybrid Search Explained: Combining BM25 and Semantic Embeddings for AI Agents
- [10] Implementing Hybrid Semantic + Lexical Search - Kent C. Dodds <https://kentcdodds.com/blog/implementing-hybrid-semantic-lexical-search>
- [11] Hybrid search explained - Redis <https://redis.io/blog/hybrid-search-explained/>
- [12] The Markdown File That Beat a \$50M Vector Database - Medium <https://medium.com/@Micheal-Lanham/the-markdown-file-that-beat-a-50m-vector-database-38e1f5113cbe>
- [13] Vector Database vs Vectorless Retrieval - Medium <https://medium.com/@vigneshkumarrajavelu/vector-database-vs-vectorless-retrieval-which-one-does-your-rag-pipeline-actually-need-f3320b989bb8>
- [14] What Is DiskANN? Billion-Scale Vector Search Explained - The Couchbase Blog <https://www.couchbase.com/blog/diskann/>
- [15] How LanceDB Accelerates Vector Search at 10 Billion Scale <https://www.lancedb.com/blog/how-lancedb-accelerates-vector-search-at-10-billion-scale>
- [16] Build billion-scale vector databases in under an hour with GPU acceleration on Amazon OpenSearch Service - AWS Big Data Blog <https://aws.amazon.com/blogs/big-data/build-billion-scale-vector-databases-in-under-an-hour-with-gpu-acceleration-on-amazon-opensearch-service/>
- [17] The Vector Database Evolution - Lawrence Emenike <https://lawrence-emenike.medium.com/the-vector-database-evolution-fe63049526b4>
- [18] From shiny object to sober reality: The vector database story, two years later - VentureBeat <https://venturebeat.com/ai/from-shiny-object-to-sober-reality-the-vector-database-story-two-years-later>
- [19] Pgvector Is Now Faster than Pinecone at 75% Less Cost - Tiger Data <https://www.tigerdata.com/blog/pgvector-is-now-as-fast-as-pinecone-at-75-less-cost>