

82%

INFERENCE COST REDUCTION

INDUSTRY

ENTERPRISE LOGISTICS

MODELS

REDIS + VERCEL AI SDK

TIMELINE

8 DAYS

STATUS

OPERATIONAL — OPTIMIZATION PHASE

Semantic Caching & API Inference Optimization

Audited a bloated internal AI tool for a logistics firm. Instituted a Redis semantic cache layer, driving API costs down by \$3,400/month.

The Baseline Inefficiency

A global logistics tech firm had rapidly deployed an internal conversational AI tool for their warehouse managers. Due to poor architectural planning and a lack of caching, identical queries regarding shipping schedules were triggering fresh API calls to the LLM every single time. Monthly API inference burn was scaling out of control, reaching \$4,100 per month with severe latency spikes during peak shift changes.

The Architectural Solution

We executed a rapid 8-day infrastructure diagnostic and deployment. We routed their Vercel AI SDK pipeline through a Redis Semantic Cache. Now, when a warehouse manager asks an identical or semantically similar query within a 12-hour window, the architecture intercepts the call and serves the cached response instantly. We also implemented Helicone for granular observability and rate-limiting.

The Fiscal Outcome

Redundant API calls were virtually eliminated. Inference costs plummeted by 82%, dropping the operational burn from \$4,100/mo to roughly \$700/mo. P99 latency was reduced to under 180ms for cached queries, stabilizing the internal tool for immediate warehouse adoption.

Quantifiable Outcomes

API OPEX REDUCTION

-82%

Drop in monthly token inference expenditure.

P99 LATENCY

180MS

Final response latency for semantically cached queries.