

ML LABS INTELLIGENCE

How to Cut AI Infrastructure Costs Safely

AI costs rise quietly through infrastructure sprawl, weak defaults, and low-visibility waste. This article shows how to cut cost without breaking reliability or performance.

TECHNICAL

Author Omar Trejo

Date 2026-02-13

ML LABS

mlabs.com/intel/cost-reduction-framework

AI infrastructure costs spiral quickly. What starts as a few hundred dollars in API calls becomes tens of thousands monthly as usage scales. Most organizations have significant cost reduction opportunities hiding in plain sight. Research from MIT confirms that **the price of a given level of benchmark performance has been decreasing by 5-10x per year** (MIT, 2025) — yet most organizations capture none of these savings because their infrastructure choices remain static.

The Hidden Cost Multipliers

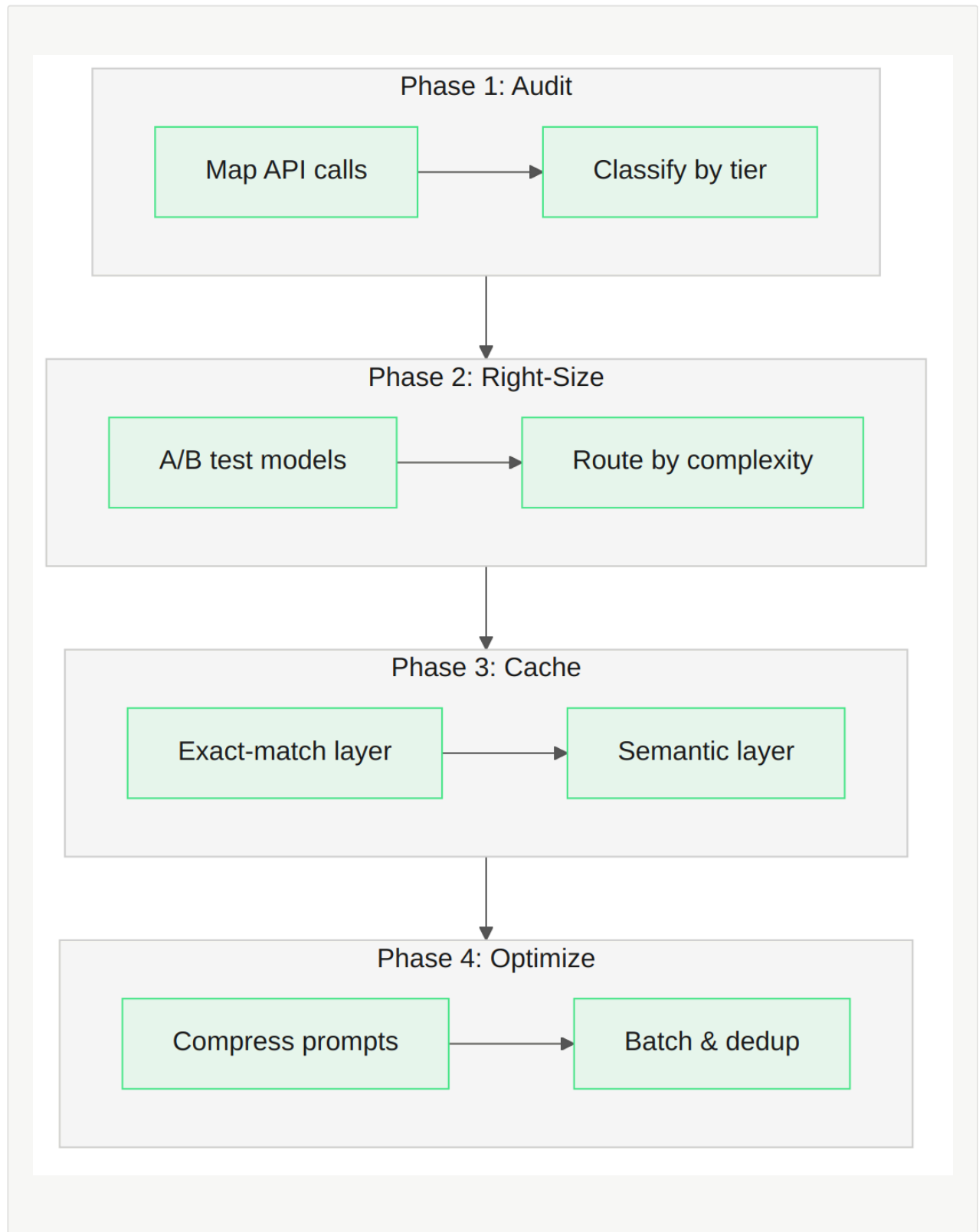
Three patterns drive the majority of unnecessary AI spend:

- **Redundant processing:** systems reprocess similar content repeatedly because they lack content-aware deduplication
- **Over-specified requests:** frontier models used for tasks that fine-tuned small models handle with equivalent accuracy
- **Inefficient prompt design:** verbose prompts inflating token costs without improving output quality

These compound as systems scale, creating cost curves that outpace business growth.

The Cost Reduction Framework

A systematic four-phase approach that builds on itself — audit first, then attack the highest-leverage areas.



Phase 1: Audit and Classify

Map every AI API call in your system: model used, average input/output tokens, frequency, and business function. Most teams discover that 20% of call patterns account for 80% of spend.

Classify each call into tiers based on quality sensitivity:

- **Tier 1 — Critical:** quality directly impacts revenue or user experience; justifies frontier models
- **Tier 2 — Important:** needs good quality but tolerates minor degradation; candidates for model downgrades
- **Tier 3 — Routine:** classification, extraction, formatting; should use the cheapest capable model

// The audit alone often reveals immediate wins. One organization discovered that 40% of their frontier model spend was on structured data extraction – a task their fine-tuned model handled identically at 1/20th the cost.

Phase 2: Right-Size Models

Run A/B tests on Tier 2 and Tier 3 calls with smaller models. Research on **LLM routing and cascading** (ICLR, 2025) demonstrates that unified model selection strategies — combining routing with cascading — consistently outperform static model assignments on cost-performance tradeoffs.

The key is measuring output quality objectively. Build evaluation harnesses that score outputs against ground truth or human ratings.

Tiered routing operationalizes this: a lightweight classifier routes each request to the cheapest model that can handle it well.

- **Classification and extraction:** small models at a fraction of frontier cost
- **Standard generation:** mid-tier models balancing quality and cost
- **Complex reasoning:** frontier models only when necessary

For high-volume, narrow tasks, **fine-tuning** a small model on your domain data often matches frontier quality at a fraction of the cost.

Phase 3: Intelligent Caching

Many queries are functionally equivalent even when phrased differently. Research on **GPT Semantic Cache** (Regmi & Pun, 2024) demonstrated that semantic caching reduced API calls by up to 69% while maintaining response quality.

Layer your caching strategy:

- **Exact match:** identical inputs return cached responses with zero compute
- **Semantic similarity:** 0.95+ cosine threshold catches paraphrases
- **Model fallback:** cache misses hit the model and populate the cache

Together, these layers eliminate up to half of all model invocations.

Phase 4: Prompt and Request Optimization

Most prompts are 2-5x longer than necessary. A **comprehensive survey of prompt compression techniques** (NAACL, 2025) found that systematic compression can achieve up to 20x reduction with minimal performance loss.

Three techniques deliver the highest impact:

- **Structured output formats:** JSON schema constraints produce identical results to verbose instructions at 10x fewer tokens
- **Content-aware deduplication:** MinHash LSH identifies near-duplicate content so only unique portions hit the model
- **Batch processing windows:** aggregating similar requests enables bulk pricing and smooths utilization

Measuring Success

Track three metrics to validate the framework is working. Cost should drop without quality degrading — if both move together, the optimization is too aggressive.

- **Cost per successful output:** total API spend divided by outputs passing quality — the number that matters most
- **Quality score distribution:** monitor the full distribution, not just the average; tail shifts signal over-optimization
- **Cache hit rate and latency:** target 35-50% cache hit rate; P50 and P99 latency should improve as a side effect

When This Approach Does Not Apply

Cost optimization without quality guardrails damages user trust. The biggest barrier is organizational, not technical — teams default to the most powerful model because it feels safer, and without measurement infrastructure to objectively evaluate output quality, cost cuts are made blind. Degradation goes undetected until stakeholders report that "the AI seems worse lately."

Before implementing cost reduction, establish quality baselines and automated monitoring for every API call tier. If building measurement infrastructure takes longer than the cost optimization itself, that is the correct sequencing.

First Steps

1. **Instrument every call.** Log cost, tokens, and quality for all AI API calls. Identify the top three cost centers by spend.
2. **Audit and test.** Classify calls by tier and A/B test your highest-volume Tier 3 calls with smaller models.
3. **Track the right metric.** Use cost-per-successful-output as the primary measure, targeting significant reduction while maintaining quality distribution.

Practical Solution Pattern

Apply cost optimization as a systematic four-phase program: audit and classify every API call, right-size model routing using A/B tests against objective quality evaluations, layer semantic caching, and compress prompts to structured output formats for the highest-volume patterns.

This works because AI costs are dominated by a small number of high-frequency, low-complexity call patterns that don't require frontier models. Right-sized routing eliminates over-specification, semantic caching removes up to half of all model invocations, and structured prompts reduce token consumption by 10x for extraction tasks with no measurable quality degradation. If you need to find where AI cost is leaking before making bigger platform changes, an **AI Operational Audit** can map the waste and prioritize the fixes.

References

1. Gundlach, H., Lynch, J., Mertens, M., & Thompson, N. **The Price of Progress: Algorithmic Efficiency and the Falling Cost of AI Inference**. MIT, 2025.
2. Dekoninck, J., Baader, M., & Vechev, M. **A Unified Approach to Routing and Cascading for LLMs**. *ICLR*, 2025.
3. Regmi, S., & Pun, C. P. **GPT Semantic Cache: Reducing LLM Costs and Latency via Semantic Embedding Caching**. arXiv, 2024.
4. Li, Z., Liu, Y., Su, Y., & Collier, N. **Prompt Compression for Large Language Models: A Survey**. *NAACL*, 2025.



ML LABS

Custom AI Systems for High-Value Workflows

mllabs.com