

The AI Cost Trap: Governance & Economic Realism

Why enterprises overpay for LLMs—and
how to fix run-rates in 30–90 days.

The Reality: You are scaling consumption before you can explain unit economics.

The Situation

Cloud spend is rising. Gartner forecasts public cloud end-user spending growth of **20.4%** in 2024. GenAI is a major contributor to this acceleration. Usage is unmetered, ungoverned, and treated as a fixed cost.

The Failure Mode

Enterprises are paying for activity, not **outcomes**. Treating AI workloads like traditional software leads to a recurring financial surprise.

The root cause is not vendor rates, but a lack of unit economic ownership.

The Fix

1. Define the unit of value.
2. Measure cost per unit.
3. Route workloads to the lowest-cost option meeting quality risk.
4. Enforce controls.

Case Study: The Cost of Ungoverned Workflows

12,000x

Overspend

Anonymized Finance Client Failure: "Training as a reflex" without economic constraints.
Result: Routine delays of days for output.



Seconds

Inference Time After Remediation

Remediation: Shifted to proper inference workflows.
Result: Slashed costs and recovered market alpha.

Takeaway: This was not a rounding error. It was a failure of operating discipline.

Diagnosis: Visible Inefficiencies in the Model Layer



1. Default-to-Premium

The Leak: Using frontier models (GPT-4) for basic classification or summarization where they aren't required.

The Fix: Define acceptance criteria and escalation rules for edge cases.



2. Training as a Reflex

The Leak: Unnecessary fine-tuning because “that's what we do,” without measuring incremental lift vs. cost.

The Fix: Require proof that inference/RAG is insufficient before approving training compute.



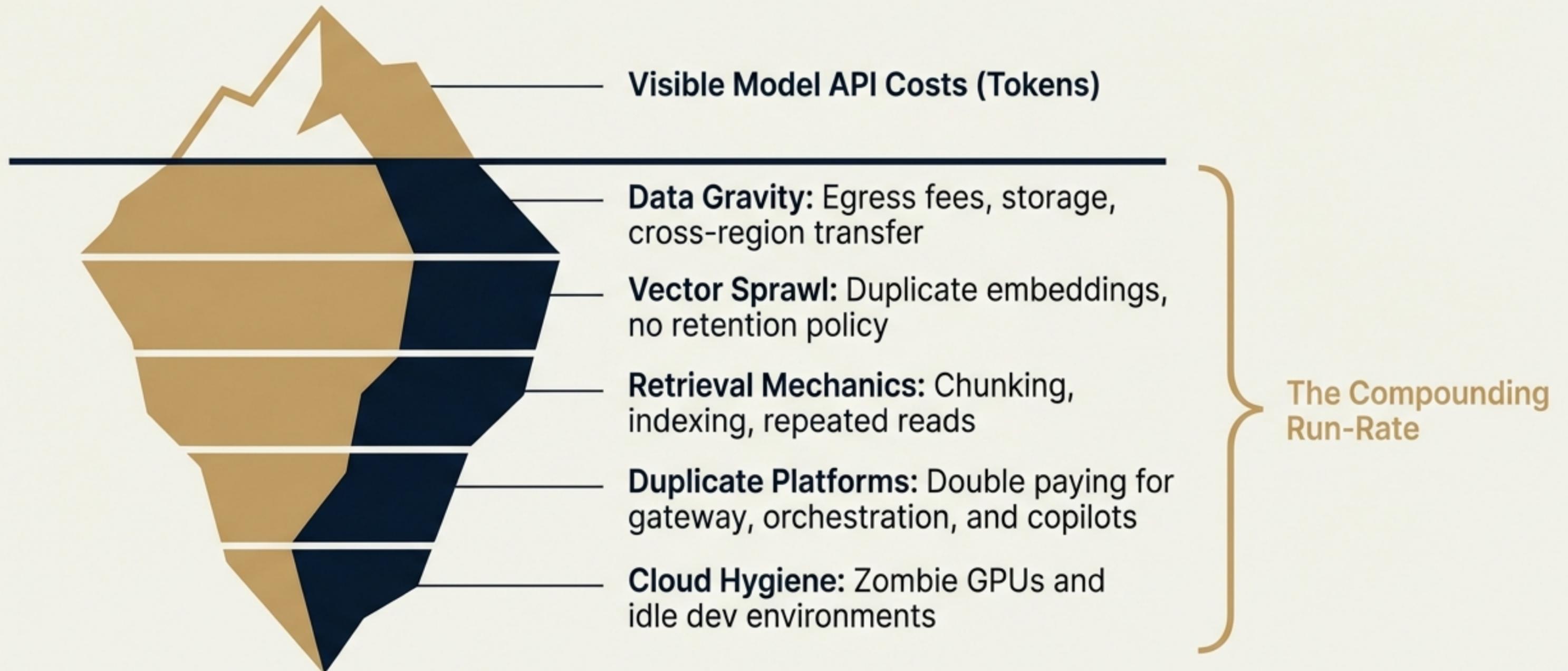
3. Token Waste

The Leak: Context bloat. Long chat histories and entire PDFs appended by default.

The Fix: Implement strict token budgets and retrieval discipline to reduce run-rate.

Diagnosis: The Hidden Multiplier (Data Gravity & Infrastructure)

Executives focus on tokens. The real cost accumulation happens below the surface.



Are you managing cost, or just observing it?

Three questions your leadership team should answer in one meeting.



1. What is our all-in cost per outcome for the top 10 workflows?

Examples: Cost per ticket deflected, cost per contract reviewed. Not just cost per token.



2. What share of usage is production vs. pilot vs. orphaned?

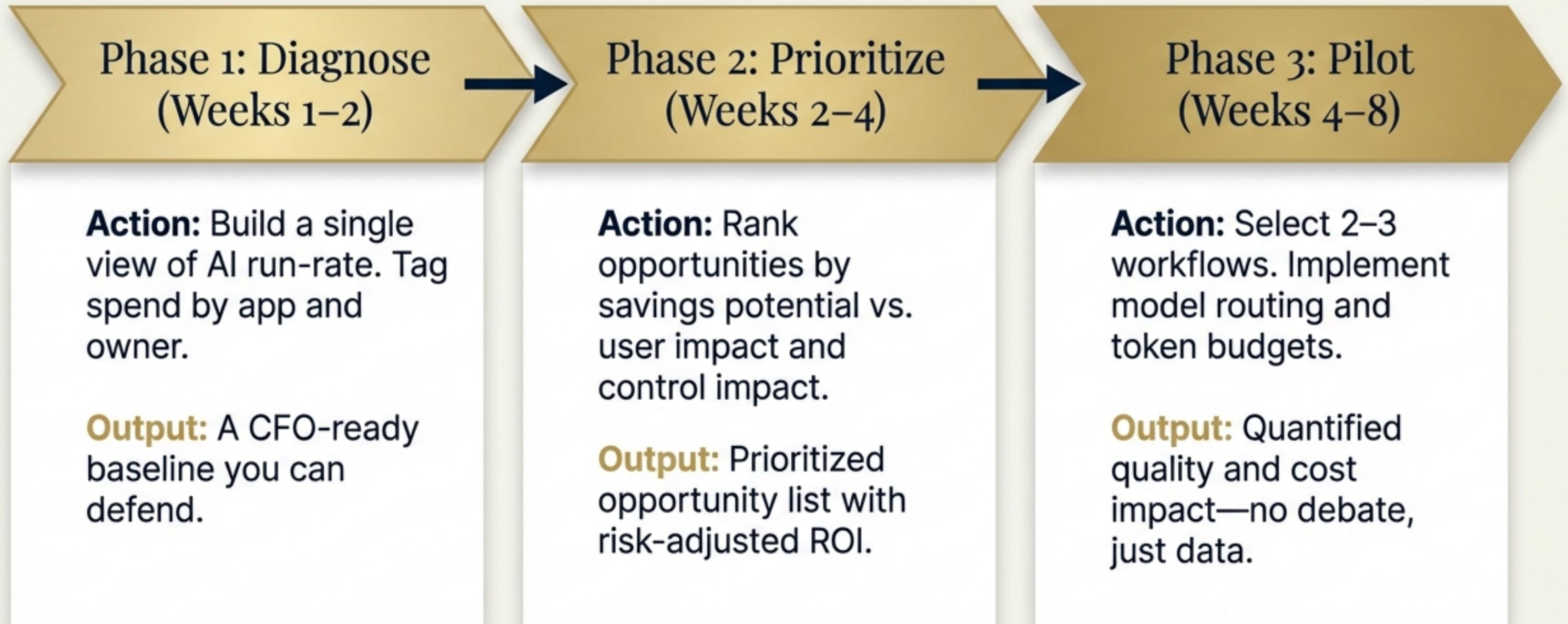
If you can't separate these, you can't cut safely.



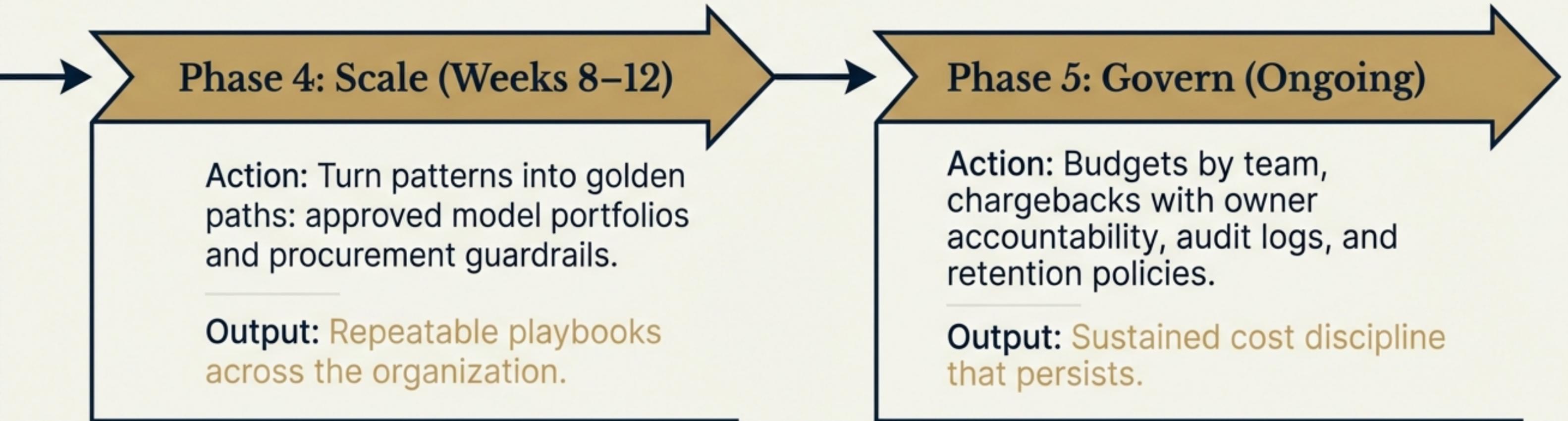
3. What is our routing policy by risk and quality tier?

Which workflows require top-tier performance, and which must default to lower-cost options?

The 30–90 Day Optimization Roadmap (Part I: Triage)



Scaling Discipline & Governance (Part II: Sustain)



PROOF POINT

In a regulated university + health system, LLM automation contributed to **\$80MM+** in cost reduction while maintaining responsible controls.

What 'Good' Looks Like: The Deliverables

Cut run-rate without cutting capability, controls, or speed.

- ✓ **AI Run-Rate Baseline:** Finance-grade reconciliation with owner tags.
- ✓ **Unit Economics Model:** Cost-per-outcome with sensitivity ranges.
- ✓ **Waste Ledger:** Quantified leaks with "stop/reduce/reneegotiate" actions.
- ✓ **Model Portfolio & Routing Policy:** Tied to risk/quality tiers with escalation rules.
- ✓ **Data Cost Controls:** Retention policies, caching strategy, and egress reduction.
- ✓ **Cloud Hygiene Plan:** Autoscaling and decommissioning rules.
- ✓ **Governance Controls:** Logging, access, and auditability standards.

Outcomes > Advice.

Independent AI value creation partner
for CEOs, CFOs, and CIOs.

Ready to find the waste?

J.L. SUTHERLAND & ASSOCIATES

contact@jlsutherland.com