

ARCHITECTURE DEEP DIVE

Trust Cascade Architecture

How intelligent routing delivers AI accuracy at rule-level cost.
The economic model for sustainable agent operations.

The Key Insight

Not every decision needs the same level of intelligence. Analysis consistently shows 60-70% of decisions can be handled by rules or simple ML. Only 5-10% genuinely require multi-agent reasoning. Route each decision to the cheapest sufficient intelligence.

65%

Handled by rules

\$0.0001

L1 cost/decision

27%

Typical cost reduction

5

Intelligence levels

IN THIS WHITEPAPER

- 01 The Agent Cost Problem
- 02 The Five Cascade Levels
- 03 Economic Model
- 04 ROI-Driven Routing
- 05 Implementation
- 06 Next Steps

01 - THE AGENT COST PROBLEM

Inference dominates everything

At 1M interactions/month, agent operations cost \$59K-\$151K. LLM inference accounts for 70-80%. Cutting observability by 50% saves \$2K. Cutting inference by 20% saves \$10-25K.

COMPONENT	MONTHLY COST	% OF TOTAL
Inference (LLM API calls)	\$45,000 - \$120,000	70-80%
Agent compute	\$8,000 - \$15,000	10-13%
Observability	\$3,000 - \$8,000	5-6%
Governance	\$2,000 - \$5,000	3-4%

02 - THE FIVE CASCADE LEVELS

Right-sized intelligence

Trust Cascade Architecture

1

Rules Engine

Deterministic rules, pattern matching, velocity checks

\$0.0001

~65%

2

ML Models

Classification, anomaly scoring, embeddings

\$0.001

~22%

3

Single Agent

LLM reasoning, tool use, structured output

\$0.02

~9%

4

Multi-Agent

Collaboration, verification, debate

\$0.08

~3%

5

Human Review

Expert escalation for high-stakes decisions

\$5.00

~1%

Confidence-Based Escalation

Each level has a confidence threshold. If a decision can be made confidently at Level 1, it stays there. If confidence is below threshold, it escalates to Level 2. And so on. Only genuinely complex decisions reach expensive levels.

03 - ECONOMIC MODEL

The math that matters

Comparing two approaches for 1 million decisions per month:

<p>ALL LLM (NO GOVERNANCE)</p> <p>\$85,000/mo</p> <p>Explicit: \$50,000 (1M x \$0.05) Hidden: \$35,000 (error remediation, compliance overhead, incident response)</p>	<p>TRUST CASCADE</p> <p>\$62,485/mo</p> <p>Explicit: \$54,485 (cascade routing) Hidden: \$8,000 (reduced remediation via governance)</p>
---	---

MONTHLY SAVINGS

\$85,000 - \$62,485 = \$22,515 (27% reduction)

04 - ROI-DRIVEN ROUTING

Decision value determines routing

Not all decisions have equal value. A \$10,000 retention decision deserves more intelligence than a \$0.10 FAQ response.

Routing Matrix: Decision Value x Complexity

	LOW VALUE (<\$10)	MEDIUM (\$10-\$1K)	HIGH (>\$1K)
LOW COMPLEXITY	Max L2	Max L3	Max L4
HIGH COMPLEXITY	Reject/Simplify	Max L4	Full Cascade + L5

Cost Ceiling Enforcement

Never spend more on a decision than it's worth. A \$0.10 decision should never route to a \$0.08 multi-agent level. Hard limits prevent economic irrationality.

Complexity Detection

High complexity + low value = product design problem. These get rejected and flagged for product team review, not thrown at expensive AI.

Dynamic Thresholds

Confidence thresholds adjust based on decision value. High-stakes decisions require higher confidence before stopping at lower levels.

Continuous Optimization

Monitor routing patterns to find optimization opportunities. If 90% of L3 decisions succeed, maybe L2 threshold is too conservative.

05 - IMPLEMENTATION WITH ROTASCALE

Trust Cascade in practice

Rotascale Orchestrate provides the infrastructure for Trust Cascade routing. Here's how each component maps:

CASCADE LEVEL	ROTASCALE COMPONENT	CAPABILITIES
L1: Rules Engine	Policy Engine	Declarative rules, pattern matching, allowlists/denylists
L2: ML Models	Guardian + Eval	Confidence scoring, anomaly detection, embedding similarity
L3: Single Agent	Orchestrate	Agent routing, tool management, structured outputs
L4: Multi-Agent	Orchestrate	Agent collaboration, verification loops, debate patterns
L5: Human Review	Human-in-the-Loop	Escalation workflows, approval queues, audit trails

Observability Built In

Every decision logged with routing path, confidence scores, cost attribution. Know exactly where your budget goes. Identify optimization opportunities.

Gradual Adoption

Start with L3/L4 (your current state). Add L1/L2 rules as patterns emerge. Shadow mode compares cascade decisions to baseline before going live.

Regulatory Compliance

Complete audit trail for every decision. Explainable routing logic. Human oversight for high-stakes decisions. Built for regulated industries.

Continuous Improvement

Analytics surface threshold tuning opportunities. Identify decisions that could move to lower levels. Track cost trends over time.

Key Insight: Governance Pays for Itself

A \$5K/month governance investment that reduces inference costs by 15% generates \$7-18K/month in savings. The Trust Cascade isn't a cost center - it's a profit center.

ROI from day one. Most deployments see positive returns within the first month.

Ready to optimize your AI economics?

Get a cost analysis for your agent operations.

[Request Analysis](#)

[Learn More](#)