FINANCIAL SERVICES

# AI That Works in Finance

*A strategic guide to deploying trustworthy AI in banking, insurance, and capital markets. From POC to production, from compliance to cost control.*

**EXECUTIVE SUMMARY**

Financial services firms are racing to deploy AI. Most will fail—not because the technology doesn't work, but because they haven't solved for trust, compliance, and economics. This guide provides a roadmap for AI that actually works in regulated environments: the frameworks, the architecture, and the operational discipline required for production success.

01 — THE CHALLENGE

# AI in Finance Is Different

Every AI decision in financial services faces scrutiny that consumer tech never sees. Explainability isn't optional—it's required. The cost of failure is measured in regulatory fines, customer trust, and market access.

> *"The question isn't whether AI can work in finance. It's whether your organization can operate AI in finance—under regulatory constraints, at enterprise scale, with the governance your stakeholders require."*

### Regulatory Burden

MAS Guidelines, OCC SR 11-7, BCBS Principles, EU AI Act. Each jurisdiction adds requirements. Each regulator expects different documentation. The compliance burden compounds.

### Explainability Requirements

Black-box decisions don't fly in finance. Customers have the right to understand decisions that affect them. Regulators have the right to audit your logic. "The model said so" isn't good enough.

### Audit Trail Expectations

Every decision documented. Every reasoning chain preserved. For years, not months. When the audit arrives, you need to reconstruct exactly what the system did and why.

### Model Risk Management

AI models are models. They need the same governance as your credit models, market risk models, and operational risk models. Continuous validation. Documented governance. MRM-compliant from day one.

## The High-Risk Categories

Under emerging AI regulations worldwide, these financial services applications are considered high-risk by default:

| APPLICATION | RISK LEVEL | KEY REQUIREMENTS |
| --- | --- | --- |
| Credit scoring & creditworthiness | High-Risk | Human oversight, explainability, bias testing |
| Insurance risk assessment & pricing | High-Risk | Documentation, validation, appeal process |
| Fraud detection (blocking transactions) | High-Risk | Accuracy monitoring, false positive review |
| Investment advice & portfolio management | High-Risk | Suitability validation, reasoning capture |
| Customer service chatbots | Medium | Transparency, handoff to human |

02 — THE POC TRAP

# Why Your Demo Worked and Production Won't

Your agentic AI POC impressed leadership. The agent answered questions, processed requests, made decisions that would have taken humans hours. Budget was approved. Six months later, you're explaining to the CFO why costs are 10x the projection.

### In a POC

- ☐ Volume is low (hundreds of requests)
- ☐ Edge cases are rare (haven't seen them yet)
- ☐ Latency is acceptable ("it's thinking!")
- ☐ Costs are invisible (lost in cloud spend)
- ☐ Monitoring is a dashboard someone checks occasionally
- ☐ Governance is "we'll figure that out later"

### In Production

- ☐ Volume is high (every request costs money)
- ☐ Edge cases are constant (and they break things)
- ☐ Latency kills UX (users won't wait 5 seconds)
- ☐ Costs are visible (and growing 20% MoM)
- ☐ Monitoring is critical (know before customers)
- ☐ Governance is mandatory (regulators are asking)

> *"The POC proved agentic AI can work. Production will prove whether you have the architecture to operate it."*

### The Five Gaps That Kill Production

**1 Cost Gap**

POC: $500/month. Production: $50,000/month. At $0.03-0.05 per decision, processing 1M transactions costs $30K-50K monthly. Most POCs never model this.

**2 Latency Gap**

In a demo, agent reasoning time is a feature. In production, it's a bug. Users expect sub-second responses. An agent making 10 LLM calls at 500ms each takes 5 seconds.

**3 Reliability Gap**

90% reliability sounds great in a POC. In production, that's 100,000 failures per million requests. And the agent doesn't throw errors—it gives confident wrong answers.

**4 Observability Gap**

The builder who knows the agent's quirks has moved on. The ops team inherits something they didn't build. When it fails at 2 AM, they have no idea why.

**5 Integration Gap**

POCs run in isolation with simplified data. Production requires real data, real security controls, real authentication. 60% of production effort is integration.

# Your AI Architecture Is Bleeding Money

Everyone focuses on cost-per-token. It's the wrong metric. The right metric is cost-per-outcome. And the 10x cost differences come from architectural decisions, not provider negotiations.

## $47K
MONTHLY LLM COST (AGENTIC)

## $3.2K
PREVIOUS SOLUTION (RULES + ML)

## 15x
COST INCREASE

Real example: A financial services client processing 1M customer interactions monthly. Their agentic AI solution cost 15x more than the previous approach. The CFO had one question: "What are we getting for that $44,000 difference?" Nobody had a good answer.

## The Math Nobody Wants to Do

| APPROACH | COST/DECISION | 1M DECISIONS | ANNUAL COST |
|---|---|---|---|
| Pure agentic (all LLM) | $0.10-0.50 | $100K-500K/mo | $1.2-6M |
| Well-architected cascade | $0.002 avg | $2K-5K/mo | $24-60K |

## The Seven Architectural Sins

### 1. Monolithic Prompts
2,000 tokens of context for every request, whether needed or not.

### 2. Retrieval Firehose
Stuffing all top-k chunks into context. Paying for tokens the model ignores.

### 3. Retry Spiral
35% of requests involve retries. That's 35% cost overhead.

### 4. Context Amnesia
No caching. Same question = same inference cost every time.

### 5. One-Model-Fits-All
GPT-4 for everything, including simple classification.

### 6. Verbose Output
Asked for yes/no, got three paragraphs. Output tokens cost 3-4x input.

**These sins multiply:** 2x (monolithic) × 1.5x (firehose) × 1.35x (retry) × 1.4x (no cache) × 1.5x (verbose) = **8.5x optimal cost**

04 — THE TRUST CASCADE

# Right-Sizing Intelligence

The solution isn't to abandon agentic AI. It's to use it where it matters. Route each decision to the cheapest sufficient intelligence layer. This is the Trust Cascade.

> "~70% of decisions can be handled by rules. ~20% by traditional ML. Only ~10% genuinely benefit from agent reasoning. But most deployments route 100% through agents. That's not strategy—that's waste."

### Level 1: Rules Engine — ~70% of decisions

Known patterns, velocity checks, policy limits. Response: <50ms. Cost: $0.0001/decision. Example: Duplicate claim detection, transaction limit checks, known bad actor lists.

### Level 2: Statistical ML — ~20% of decisions

Anomaly detection, risk scoring, pattern recognition. Response: <500ms. Cost: $0.001/decision. Example: Unusual billing patterns, geographic anomalies, network analysis.

### Level 3: Single Agent — ~7% of decisions

Complex reasoning, document analysis, policy interpretation. Response: 2-3s. Cost: $0.01/decision. Example: Medical necessity evaluation, multi-document synthesis.

### Level 4: Multi-Agent Tribunal — ~3% of decisions

Adversarial debate for high-stakes, ambiguous cases. Response: 3-5s. Cost: $0.03-0.05/decision. Example: Complex fraud cases requiring prosecution, defense, and judge agents.

## The Economics Transform

| LAYER | VOLUME | COST/DECISION | MONTHLY COST (1M) |
|---|---|---|---|
| L1: Rules | 700,000 | $0.0001 | $70 |
| L2: ML | 200,000 | $0.001 | $200 |
| L3: Single Agent | 70,000 | $0.01 | $700 |
| L4: Multi-Agent | 30,000 | $0.04 | $1,200 |
| Total | 1,000,000 | $0.0022 avg | $2,170 |

**Result:** Same 1M decisions. $2,170/month instead of $100,000+. 98% cost reduction with no accuracy loss.

## 94%
DETECTION ACCURACY MAINTAINED

## 98%
COST REDUCTION VS PURE AGENTIC

## 16 wk
IMPLEMENTATION TIMELINE

# The Global Regulatory Wave

AI regulation isn't coming—it's here. From the EU to Singapore to the US, regulators are converging on similar requirements: governance, transparency, human oversight. The question isn't whether you'll need to comply, but whether you're building for it now.

| MAS Guidelines | OCC SR 11-7 | BCBS Principles | EU AI Act |
|---|---|---|---|
| Singapore's responsible AI for financial services | US model risk management requirements | Basel Committee AI governance standards | Comprehensive high-risk AI regulation |

## Where Regulation Stands Today

**Europe**

**EU AI Act (2024-2026)**

Most comprehensive framework. High-risk AI systems in financial services face strict requirements by August 2026.

**United States**

**OCC SR 11-7 + State Laws**

Model risk management for banks. Colorado, California, and others adding AI-specific rules.

**Asia-Pacific**

**MAS, HKMA, RBI Guidelines**

Singapore, Hong Kong, and India leading with principle-based AI governance for financial institutions.

## What High-Risk Requires

**Risk Management**

Ongoing process to identify, estimate, and mitigate risks. Documented and updated throughout lifecycle.

**Data Governance**

Training data must be relevant, representative, error-free. Document sources, processing, assumptions.

**Technical Documentation**

Purpose, function, data, testing results. Detailed enough for regulatory evaluation.

**Record-Keeping**

Automatic logging of operations. Reconstruct what happened and why. Years of retention.

**Transparency**

Users know they're interacting with AI. Deployers can explain outputs to affected individuals.

**Human Oversight**

Effective human oversight capability. Understand outputs, decide action, override when needed.

**The good news:** Regulators worldwide are converging on the same principles—risk management, documentation, human oversight. Build governance once, and you're positioned for compliance everywhere.

# The Five Gates

Before any AI feature launches in financial services, it must clear five gates. These aren't bureaucratic hurdles—they're the foundations of trustworthy AI operations.

**1**    **Gate 1: Evaluation**

Do you know if it works? Test dataset exists. Baseline metrics documented. Automated eval suite runs in CI/CD. Deployments gated on eval pass. Edge cases and adversarial inputs tested.

**2**    **Gate 2: Monitoring**

Will you know when it breaks? Latency, errors, tokens tracked. Quality scoring in production. Drift detection enabled. Alerts configured with on-call routing.

**3**    **Gate 3: Fallbacks**

What happens when it fails? Timeouts and retries configured. Backup model available. Graceful degradation defined. Human escalation path exists.

**4**    **Gate 4: Cost Controls**

Can you afford it at scale? Budget and spending cap defined. Rate limits implemented. Cost tracking enabled. Alerts at 50/75/90% of budget.

**5**    **Gate 5: Operations**

Can you respond to incidents? Runbooks written. On-call rotation established. Rollback procedure tested. Post-mortem process defined.

> *"A feature should not launch until all five gates are cleared. This isn't optional in financial services—it's the minimum bar for regulatory compliance and operational reliability."*

## The Fallback Hierarchy

| IF THIS FAILS... | THEN TRY THIS | LAST RESORT |
| --- | --- | --- |
| Primary model | Retry with backoff | Backup model |
| Backup model | Check semantic cache | Return cached response |
| No cache hit | Human escalation | Static fallback response |

# The Shadow AI Crisis in Banking

Every major bank is deploying AI agents. Trading desks have them. Operations has them. Compliance has them. Customer service has them. But who's governing them?

> *"We had 47 AI agents running across the organization. Nobody knew what they were doing, who owned them, or what decisions they were making. Regulators started asking questions we couldn't answer."*
>
> — Chief Risk Officer, Global Systemically Important Bank

## The Ungoverned Reality

This is what we see at every major financial institution:

### Shadow Agents Everywhere

Business units deploy AI agents without IT approval. Data science teams experiment in production. Vendors embed agents in software. Nobody has a complete inventory of what's running.

### No Decision Visibility

Agents make thousands of decisions daily. Trading recommendations. Customer responses. Risk assessments. When something goes wrong, nobody can reconstruct what happened or why.

### Policy Violations After the Fact

Compliance reviews agent outputs manually—days or weeks later. By then, the customer got bad advice. The trade executed. The damage is done. Reactive, not preventive.

### Untracked Spending

API costs scattered across departmental budgets. No attribution by use case. CFO asks "what are we spending on AI?" Nobody has a real answer. Millions in hidden costs.

## What Regulators Are Asking

**"Show us your AI inventory"**
Complete list of all AI systems, their purposes, and risk classifications

**"Explain this decision"**
Full reasoning chain for any customer-impacting AI output

**"Prove your controls"**
Evidence that policies are enforced, not just documented

**"Who approved this?"**
Human oversight evidence for high-risk operations

**The question isn't whether you need agent governance. It's whether you'll implement it before or after the regulatory inquiry.**

AGENTOPS

# Enterprise Agent Governance Platform

AgentOps transforms shadow AI chaos into governed, auditable, compliant operations. Every agent registered. Every decision recorded. Every policy enforced.

## The Transformation

| BEFORE AGENTOPS | AFTER AGENTOPS |
|---|---|
| Shadow agents deployed without oversight | Every agent registered with URN-based identity |
| No visibility into agent decisions | Flight Recorder captures every reasoning step |
| Policy violations discovered after the fact | Three-layer policy enforcement in real-time |
| Untracked AI spending across departments | Full cost attribution by agent and use case |
| No audit trail for regulators | Complete decision lineage for compliance |
| Agent dependencies unknown | Composability mapping shows all relationships |

## Core Capabilities

### Agent Registry & Identity

Central inventory of every AI agent in your organization. URN-based identification (org/domain/name/version) tracks agents across dev, staging, and production. Full lineage from creation through retirement. Filter by status, autonomy level, risk classification, or owning department.

### Flight Recorder

Step-by-step visualization of every agent decision. See the complete chain: input data → context retrieval → policy evaluation → reasoning steps → human approval (if required) → final output. Replay any decision for audit. Export for regulatory inquiry. Litigation-ready documentation.

### Policy Enforcement

Three enforcement layers that work together: Gateway (block requests before they reach agents), Sidecar (monitor and intervene at runtime), Inline (enforce constraints within agent context). Human-in-loop workflows for high-risk operations. Real-time compliance scoring.

### Composability Mapping

Visual graph of agent-to-agent communication patterns. See which agents call other agents. Understand data flows and dependencies. Clustering analysis reveals hidden relationships. Critical for understanding systemic risk in your agent ecosystem.

### See AgentOps in Action

Interactive demo with real governance scenarios

**Launch Demo**

# Trust Intelligence for Finance

The Rotascale Trust Intelligence Platform provides the complete infrastructure for trustworthy AI in financial services. Seven products. One integrated platform.

### Guardian

AI reliability monitoring. Sandbagging detection, hallucination monitoring, drift detection. 96% detection accuracy. Know when your AI is failing before customers do.

### Steer

Runtime behavior control. Adjust model outputs without retraining using steering vectors. Enforce compliance language. Prevent policy violations in real-time.

### Eval

LLM evaluation platform. Rigorous, reproducible testing at scale. CI/CD integration. MRM-compliant test documentation. Automated regression detection.

### Orchestrate

Multi-agent platform with Trust Cascade. Visual builder for agent workflows. Verification protocols. Complete audit logs. Cost-optimized routing.

### Context Engine

Context-first data processing. ETL-C pipelines that preserve meaning for AI consumption. Financial data connectors. Real-time context retrieval.

### AgentOps

Enterprise agent governance. Registry, lifecycle management, policy enforcement, flight recorder. The control plane for AI agents at scale.

### Accelerate

Inference optimization. Speculative decoding for 2-8x speedup on production workloads. Lower latency, lower costs, same accuracy.

## Built for Financial Services

### Your Infrastructure

On-premise, private cloud, or hybrid deployment. No data leaves your environment. Air-gapped deployment available for sensitive workloads.

### Compliance-Ready

SOC 2 compliant. Audit logs designed for global AI regulations. MRM-compatible documentation. Built for regulatory scrutiny.

### Research-Backed

Every product built on peer-reviewed methods from Rotalabs. Open source foundations. Verify our claims at rotalabs.ai.

USE CASES

# What Banks Build With the Platform

Real implementations across trading, operations, compliance, and customer service.

### Fraud, Waste & Abuse Detection

Trust Cascade routes each transaction to the cheapest sufficient layer. Rules catch 70% instantly. ML handles 20% in milliseconds. Agents reason on the 10% that matter.

**Result:** 94% detection at 86% lower cost than pure agentic approaches.

### AI-Assisted Underwriting

Context Engine synthesizes data from dozens of sources—credit bureaus, financial statements, market data. Orchestrate coordinates multi-agent analysis. Full reasoning capture for adverse action documentation.

**Result:** 60% faster decisions with complete audit trail.

### Customer Service Concierge

24/7 AI assistant that knows account details, answers questions instantly, handles service requests. Guardian monitors for hallucination in real-time. Steer enforces compliance language.

**Result:** 40% call deflection with higher CSAT scores.

### Transaction Trust

Multi-tier decision routing for transaction approval. Simple transactions auto-approved by rules. Complex transactions escalate through ML and agent layers. Complete reasoning chains preserved.

**Result:** Sub-second approval for 95% of transactions.

### Model Risk Management

Continuous evaluation for MRM compliance. Eval runs pre-deployment testing at scale. Guardian tracks runtime drift and confidence calibration. Compliance dashboard for SR 11-7, MAS, and global AI regulations.

**Result:** Automated MRM documentation, continuous validation.

### Investment Advisory

Robo-advisory with human oversight checkpoints. Suitability validation built in. Reasoning capture for every recommendation. Steer ensures advice stays within approved parameters.

**Result:** Scalable advice with full explainability.

> "The platform doesn't replace your AI strategy—it makes your AI strategy production-ready. Same capabilities, but with the trust infrastructure financial services requires."

# Engagement Options

Every financial services AI journey is different. We offer engagements at every stage—from strategy through production, from assessment through ongoing advisory.

### AI Readiness Assessment

## $15K

2 weeks. Current state audit. Gap analysis against production readiness. Prioritized roadmap. Business case development.

### AI Strategy Workshop

## $25K

3 weeks. Use case prioritization. Architecture design. Cost modeling. Compliance alignment. Executive presentation.

### Production Implementation

## $30K+

8-16 weeks. Full platform deployment. Integration with your systems. Team training. Go-live support.

## For Specific Use Cases

| ENGAGEMENT | DURATION | INVESTMENT | DELIVERABLE |
|---|---|---|---|
| FWA Detection Assessment | 2-3 weeks | $30K | Detection audit, cascade design, business case |
| FWA Pilot Implementation | 6-8 weeks | $75K | Working cascade for one claim type |
| FWA Production Platform | 4-6 months | $300K+ | Complete cascade, integrations, training |
| Ongoing Advisory | Monthly | $8K/mo | Architecture reviews, optimization, support |

### Ready to build AI that works in finance?

Let's discuss your specific requirements and challenges.

Request Consultation

See Demos

### Contact

contact@rotascale.com · +1 (415) 524-0007

rotascale.com/solutions/financial-services