

Securing Enterprise AI: Governance, Risk, and Economics of LLM Adoption

By : Murali Jagdev Koney, Technical Team Lead at Cox Automotive

Conf42.com DevSecOps 2025

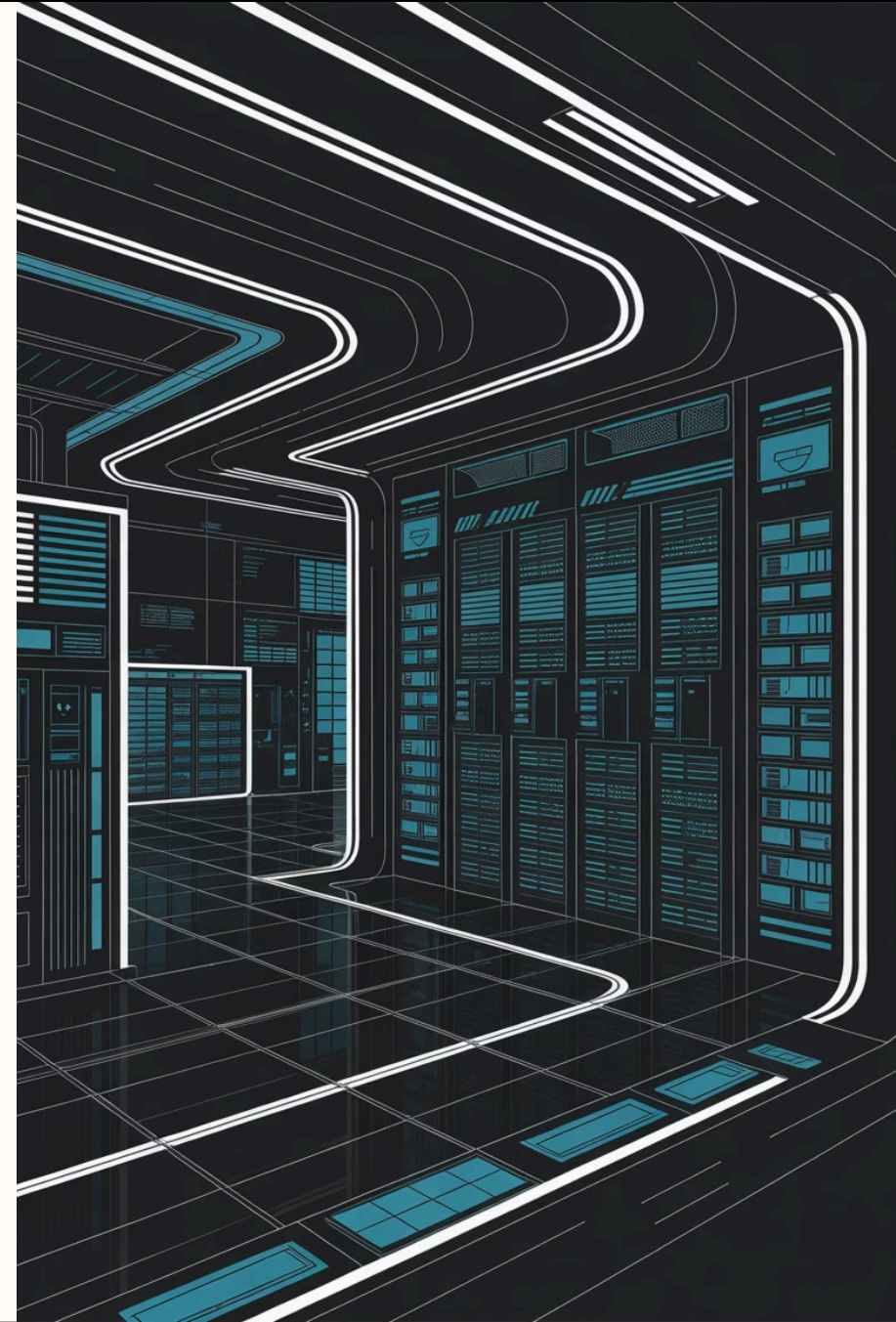
The Enterprise AI Reality Check

The Promise

Large Language Models are transforming how enterprises operate, promising enhanced productivity, intelligent automation, and competitive advantage. Organizations across industries are racing to integrate AI capabilities into their core operations.

The Challenge

Most LLM initiatives fail to reach production. The gap between proof-of-concept and production deployment remains wide, with infrastructure complexity, governance gaps, and security misalignment creating formidable barriers for DevSecOps teams.



Why Most LLM Initiatives Fail

Infrastructure Gaps

Lack of scalable, secure infrastructure designed for AI workloads. Traditional systems struggle with the compute and storage demands of LLM operations.

Governance Vacuum

Absence of clear policies for AI model selection, data handling, and decision-making accountability. Compliance frameworks lag behind AI capabilities.

Security Misalignment

Security teams unprepared for AI-specific threats like prompt injection, model poisoning, and data exfiltration through model outputs.

Cost Uncertainty

Unpredictable expenses and unclear ROI models make it difficult to justify continued investment and scale beyond initial pilots.



The Three-Pillar Framework for Secure LLM Deployment

Strategic Infrastructure

Choosing the right deployment model: IaaS, PaaS, or SaaS. Each offers distinct trade-offs between control, compliance requirements, and speed to market.

Economic Analysis

Rigorous TCO modeling that accounts for compute costs, licensing, operational overhead, and performance optimization across the AI lifecycle.

Operational Safeguards

LLMOps practices and gateway architectures that enforce security policies, monitor behavior, and prevent costly failures in production.

Deployment Models: IaaS vs PaaS vs SaaS

Infrastructure as a Service

Maximum Control: Full control over model selection, fine-tuning, and infrastructure configuration. Ideal for highly regulated environments with specific compliance needs.

Trade-offs: Highest operational burden, longer time-to-market, requires deep ML expertise and significant DevOps investment.

Platform as a Service

Balanced Approach: Managed infrastructure with flexibility for customization. Faster deployment than IaaS while maintaining significant control over model behavior.

Trade-offs: Some vendor lock-in, limited infrastructure customization, but dramatically reduced operational complexity.

Software as a Service

Fastest Time-to-Market: Pre-built solutions with minimal setup. Perfect for standard use cases and rapid experimentation with minimal technical overhead.

Trade-offs: Least control over model behavior, potential compliance concerns, limited customization for domain-specific needs.

Five-Year TCO Analysis Framework

Total Cost of Ownership extends far beyond initial licensing or compute costs. A comprehensive five-year analysis must account for infrastructure, operational overhead, talent acquisition, and ongoing optimization efforts.

Understanding how utilization rates and performance metrics directly impact your bottom line is critical for making informed deployment decisions and setting realistic ROI expectations.

01

Infrastructure Costs

Compute resources, storage, networking, and scaling requirements

02

Licensing & Models

API costs, model licensing, or hosting fees

03

Operational Overhead

Monitoring, maintenance, security, and compliance

04

Talent & Training

Specialized AI/ML engineers and ongoing education

05

Optimization Efforts

Continuous improvement, fine-tuning, and efficiency gains

Performance Metrics That Drive ROI

First-Token Latency

Time until the model begins responding. Critical for user experience in interactive applications. Lower latency drives higher adoption and satisfaction.

Tokens Per Second

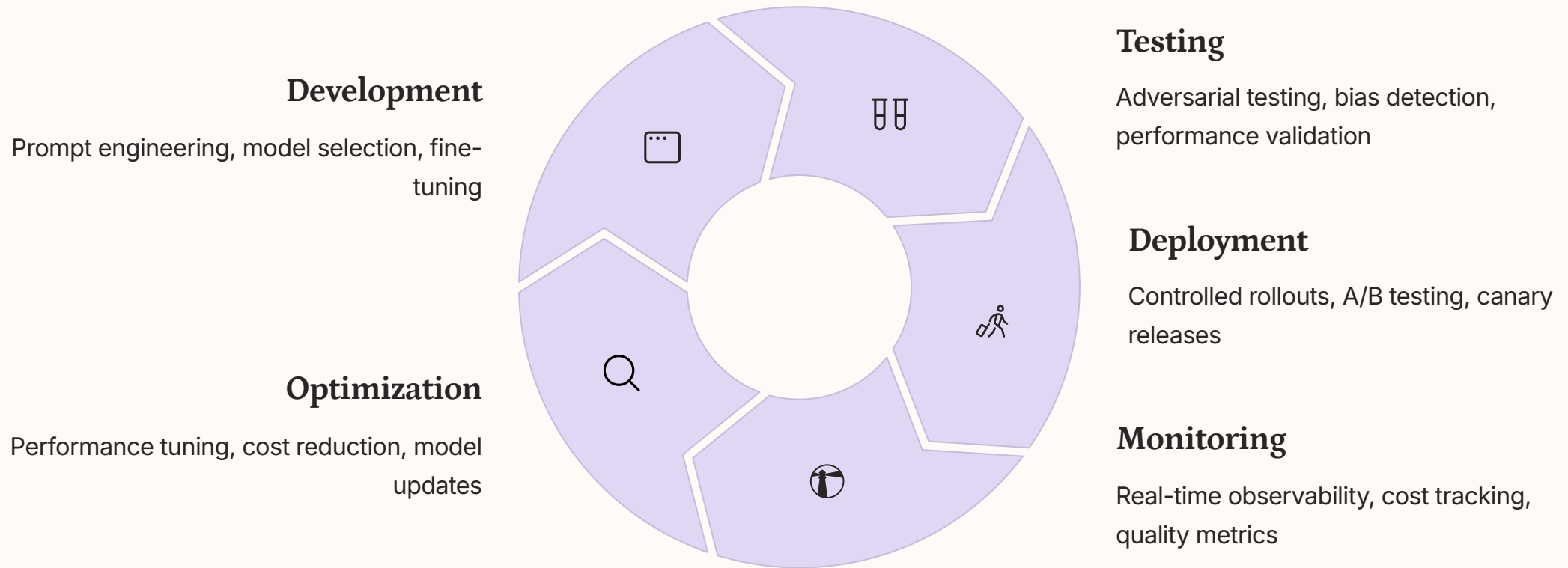
Throughput capacity determines how many concurrent users you can serve. Higher throughput reduces infrastructure needs and improves cost efficiency.

Utilization Rate

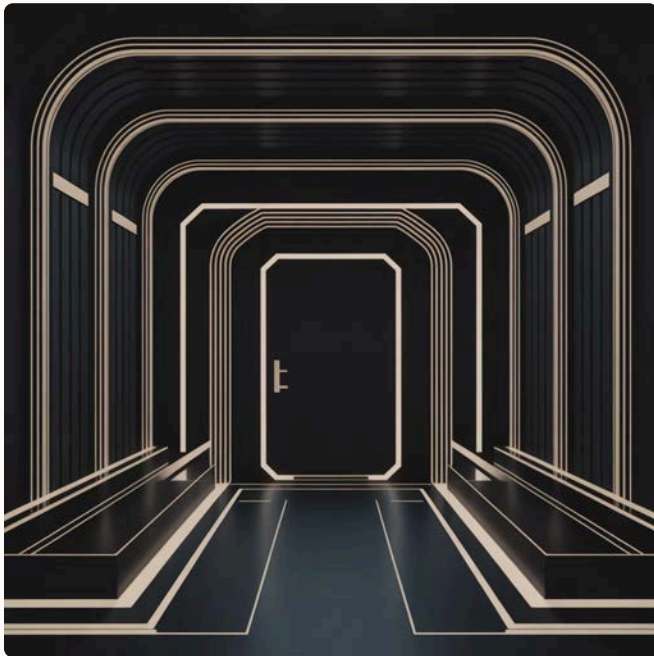
Percentage of provisioned capacity actively serving requests. Optimizing utilization is the fastest path to improving ROI on LLM infrastructure investments.

LLMOps: Operational Excellence for AI

LLMOps extends traditional MLOps practices to address the unique challenges of Large Language Models. It encompasses the end-to-end lifecycle: from prompt engineering and model selection through deployment, monitoring, and continuous improvement.



AI API Gateways: Your Security Control Plane



Essential Gateway Capabilities

AI API Gateways serve as the critical control point for all LLM interactions, enforcing security policies and preventing abuse before requests reach your models.

→ Authentication & Authorization

Verify identity and enforce role-based access control

→ Rate Limiting & Quotas

Token-based throttling prevents cost overruns and abuse

→ Input Sanitization

Detect and block prompt injection attempts

→ Output Filtering

Scan responses for sensitive data before returning

→ Semantic Caching

Reduce costs by caching similar queries

Critical Security Risks in LLM Operations

1

Prompt Injection Attacks

Malicious users craft inputs that manipulate model behavior, bypass safety controls, or extract training data. Gateway-level filtering and input validation are essential defenses.

2

Data Exfiltration

Models may inadvertently expose sensitive information from training data or context windows. Output filtering and access controls limit exposure.

3

Hallucination Risks

Models generate plausible but incorrect information, creating liability in high-stakes domains. Confidence scoring and human-in-the-loop validation mitigate this risk.

4

Cost Overruns

Uncontrolled usage can rapidly escalate costs. Token-based rate limiting and budget alerts prevent financial surprises.

Domain-Specific Models for Regulated Industries

Generic models often fall short in regulated sectors where accuracy, compliance, and domain expertise are non-negotiable. Domain-specific models trained on industry data provide superior performance and built-in compliance advantages.



BloombergGPT

Financial services model trained on decades of market data, earnings calls, and financial documents. Delivers accurate analysis while understanding financial terminology and regulatory context.



Med-PaLM

Healthcare model achieving expert-level performance on medical licensing exams. Supports clinical decision-making while maintaining HIPAA compliance and medical accuracy standards.



Best Practices for Production LLM Deployment

- **Implement Adversarial Testing**

Regularly test models with malicious inputs and edge cases to identify vulnerabilities before attackers do

- **Deploy Semantic Caching**

Cache responses to similar queries, reducing costs and improving response times without sacrificing quality

- **Enforce Token-Based Rate Limiting**

Control costs and prevent abuse by limiting tokens per user, per time period, with intelligent quota management

- **Establish Human-in-the-Loop Workflows**

Critical decisions require human oversight, particularly in regulated industries or high-stakes scenarios

- **Monitor Continuously**

Track performance metrics, cost trends, and security events in real-time with automated alerting

- **Version Control Everything**

Treat prompts, model configurations, and policies as code with proper versioning and rollback capabilities

- **Implement Fallback Strategies**

Design graceful degradation when models are unavailable or produce low-confidence outputs

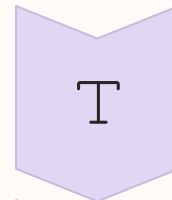
- **Document Decision Trails**

Maintain audit logs of model decisions for compliance, debugging, and continuous improvement



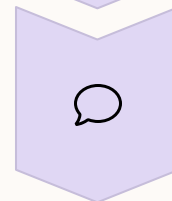
The Future: Multimodal and Agentic AI

The next evolution of enterprise AI moves beyond text-only models to multimodal systems that process images, audio, video, and text simultaneously. Agentic AI systems can autonomously plan, execute tasks, and interact with external tools.



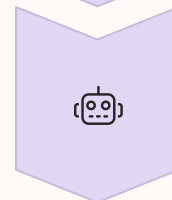
Text-Only LLMs

Current generation: powerful but limited to written language



Multimodal Models

Process multiple data types simultaneously for richer understanding



Agentic Systems

Autonomous agents that plan, execute, and adapt to achieve goals

DevSecOps teams must prepare for these advances by building flexible, secure infrastructure that can accommodate rapid AI evolution while maintaining governance and compliance standards.

From Tools to Ecosystems: The Integration Imperative

The Paradigm Shift

Successful enterprises are moving beyond treating LLMs as isolated tools and instead integrating them as governed ecosystems within their DevSecOps pipelines.

This means embedding AI capabilities directly into development workflows, security scanning, incident response, and operational monitoring—with consistent governance and observability across all touchpoints.

Key Integration Points

- CI/CD pipelines for automated code review and security scanning
- Incident response systems for intelligent triage and remediation
- Monitoring platforms for anomaly detection and root cause analysis
- Documentation systems for automated knowledge base generation



Thank You!

Murali Jagdev Koney

Technical Team Lead, Cox Automotive

Conf42.com DevSecOps 2025