

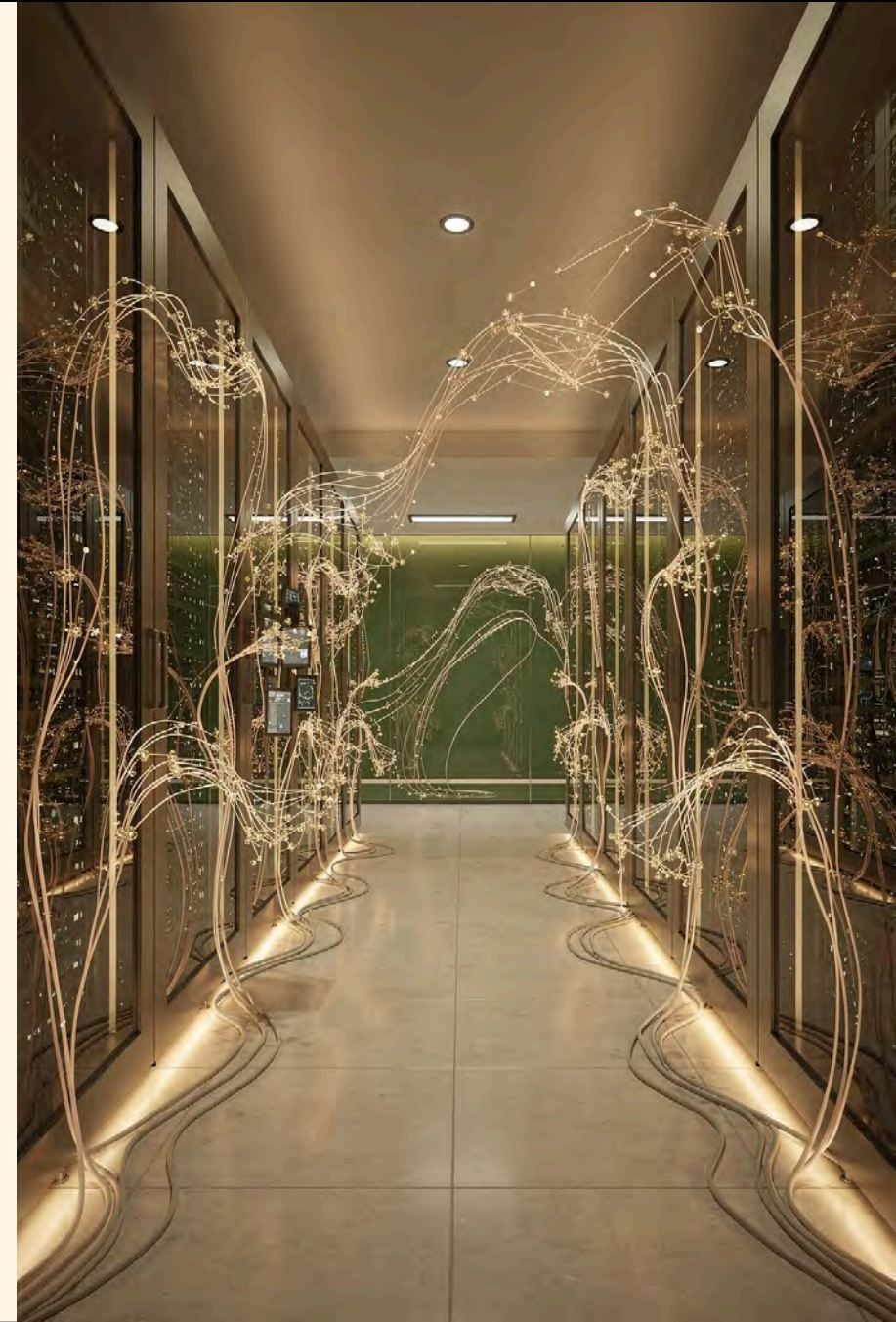
From Reliable Models to Resilient ML Platforms: Designing secure, highly available infrastructure for production machine learning systems



CONF MACHINE LEARNING

SRINIVAS TALASILA

SAP AMERICA INC



Today's Agenda

01

The Production ML Challenge

Understanding the complexity of production ML systems

02

Modernizing Infrastructure

Moving from legacy to cloud-native architectures

03

IBM Cloud SoftLayer Foundation

Building on enterprise-grade infrastructure

04

Pillars of Resilience

High availability, security, and performance

05

People & Process

Cross-functional collaboration and responsibility models

06

Deployment Patterns

Cloud-native and hybrid approaches

07

Key Takeaways

Essential lessons for resilient ML platforms



The Production ML Challenge

Moving machine learning from experimentation to production introduces critical infrastructure demands. Model performance alone isn't enough systems must be reliable, available, and trustworthy at scale. Legacy infrastructure often fails to meet the performance and security requirements of continuous ML workloads like training, batch inference, and real-time prediction services.

Modernizing Infrastructure for ML at Scale

Legacy Systems

Rigid, monolithic infrastructure creates critical bottlenecks, leading to single points of failure and an inherent inability to scale elastically. These systems struggle to process vast datasets efficiently, severely hindering distributed ML workloads.

Production Demands

The complexity of operationalizing ML models necessitates sophisticated infrastructure. This includes managing intricate continuous training pipelines, executing real-time inference at sub-millisecond latencies, supporting robust model versioning and A/B testing at scale, integrating feature stores, proactively monitoring for model drift, and gracefully handling unpredictable peak loads.

Modern Platform

A cutting-edge platform delivers a resilient, fault-tolerant architecture designed for elastic scaling. It provides comprehensive observability for proactive issue detection, enables automated recovery from failures, and ensures seamless CI/CD integration for rapid, reliable model deployment across distributed systems.

The Foundation: IBM Cloud SoftLayer



Compute Layer Strategy

Leveraging IBM Cloud SoftLayer as the foundational compute layer enabled seamless integration across distributed ML systems. The platform provides the performance baseline required for data-intensive pipelines while supporting flexible scaling patterns.

- High-performance bare metal and virtual compute options
- Global availability zones for distributed training
- Network backbone optimized for large data transfers

Pillars of Resilient ML Infrastructure



High Availability

Multi-zone deployment patterns ensure ML workloads continue uninterrupted during infrastructure failures or maintenance windows



Disaster Recovery

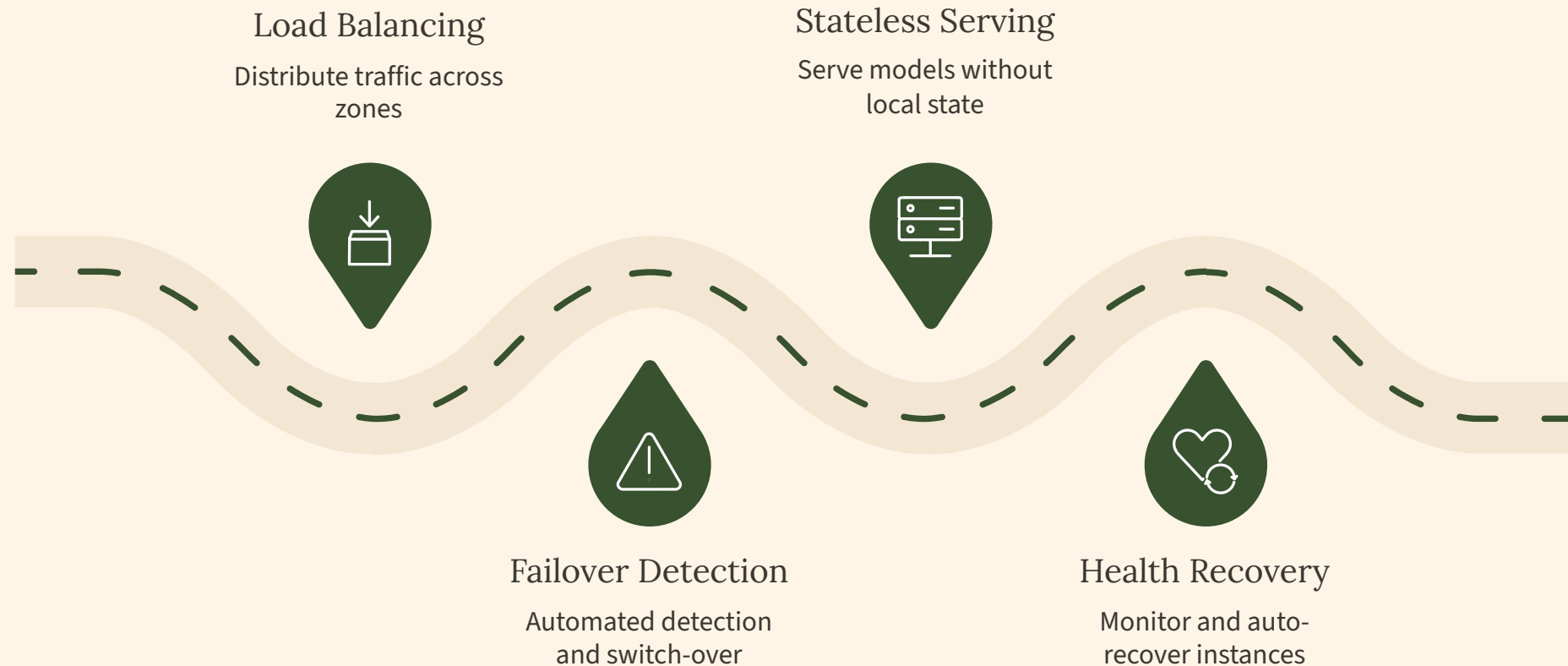
Automated backup strategies and cross-region replication protect model artifacts, training data, and configuration state



Security by Design

Zero Trust principles embedded from inception isolate environments and protect sensitive ML assets throughout the pipeline

High Availability Architecture



Distributed deployment across availability zones eliminates single points of failure in model serving infrastructure.

Sustaining ML Workloads

Production ML demands continuous availability across various critical workload types, each presenting unique challenges for robust system design:

- **Model Training:** Long-running distributed jobs across GPU clusters requiring fault-tolerant checkpointing, node failure recovery, and coordinated parameter updates.
- **Batch Inference:** Scheduled pipelines processing millions of records with complex data dependencies, variable volumes, and strict downstream SLA requirements.
- **Real-time Serving:** Low-latency API endpoints with sub-second response times, handling unpredictable traffic spikes through dynamic auto-scaling across geographic regions.



Zero Trust Segmentation for ML

Security embedded at the infrastructure layer protects ML systems from lateral movement and unauthorized access. Zero Trust Segmentation creates isolated zones for training, inference, and data storage each with explicit access controls and continuous verification.

Training Environment Isolation

Separate compute clusters prevent training workloads from accessing production serving infrastructure

Dataset Protection

Sensitive training data encrypted at rest and in transit, with access limited by role and workload identity

Model Serving Security

Inference endpoints protected by API authentication, rate limiting, and network segmentation policies

Cloud Security Framework Alignment

Operational Controls

Aligning cloud security frameworks with ML operational requirements protects assets while maintaining experimentation velocity.

- Identity and access management for model registries
- Audit logging across training and serving pipelines
- Compliance validation for regulated ML workloads
- Automated security scanning of container images



PERFORMANCE METRICS

Measuring Platform Resilience

Model Serving Uptime

Availability across production
inference endpoints

Failover Recovery

Automated recovery time for zone-
level failures

Training Throughput

Improvement in distributed training
performance

Security Breaches

Lateral movement incidents since segmentation deployment



Beyond Infrastructure: People and Process

Resilient ML operations require more than technology they demand a cultural shift toward shared responsibility across organizational boundaries.

When security and reliability ownership extends beyond infrastructure teams to include data scientists, ML engineers, MLOps practitioners, and platform engineers, organizations create systems where resilience is built-in at every stage from model development to deployment and monitoring.

This distributed ownership model ensures teams collaborate on shared outcomes rather than simply handing off responsibilities. The result: ML systems that are robust, secure, and scalable by design, not as an afterthought.

Cross-Functional Responsibility Model

Production ML reliability emerges from clear accountability across specialized teams working in concert.

Data Science

Model monitoring, performance validation, retraining triggers

MLOps

Pipeline reliability, deployment automation, feature store management

Platform Engineering

Infrastructure HA/DR, security controls, capacity planning



Deployment Patterns: Cloud-Native and Hybrid

Platform Flexibility

The architectural principles demonstrated apply across deployment models. Whether operating in fully cloud-native environments or managing hybrid infrastructure, the same HA, DR, and security patterns ensure ML platform resilience.

Hybrid deployments introduce complexity around data gravity, latency requirements, and compliance boundaries but the fundamental design principles remain consistent.

Practical Lessons for Production ML

Infrastructure Foundations Matter

Choose compute platforms that support both performance and resilience requirements from the start

Embed Security Early

Retrofit security is expensive and risky design Zero Trust segmentation into initial architecture

Design for Failure

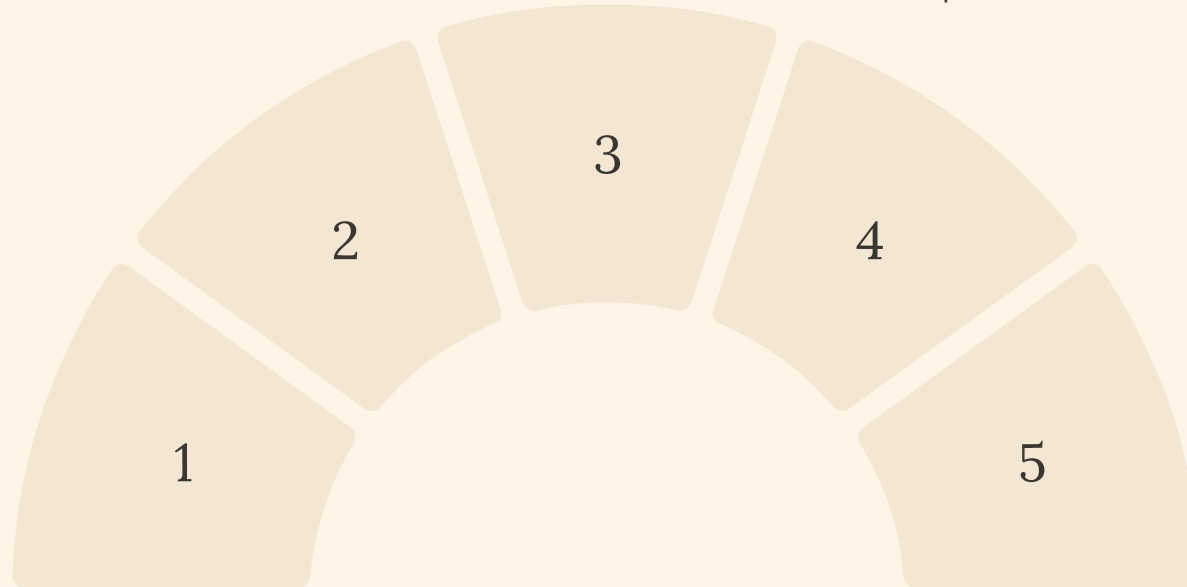
Assume infrastructure components will fail and architect ML systems to handle failures gracefully

Distribute Responsibility

Reliability emerges from shared ownership across data science, MLOps, and platform teams

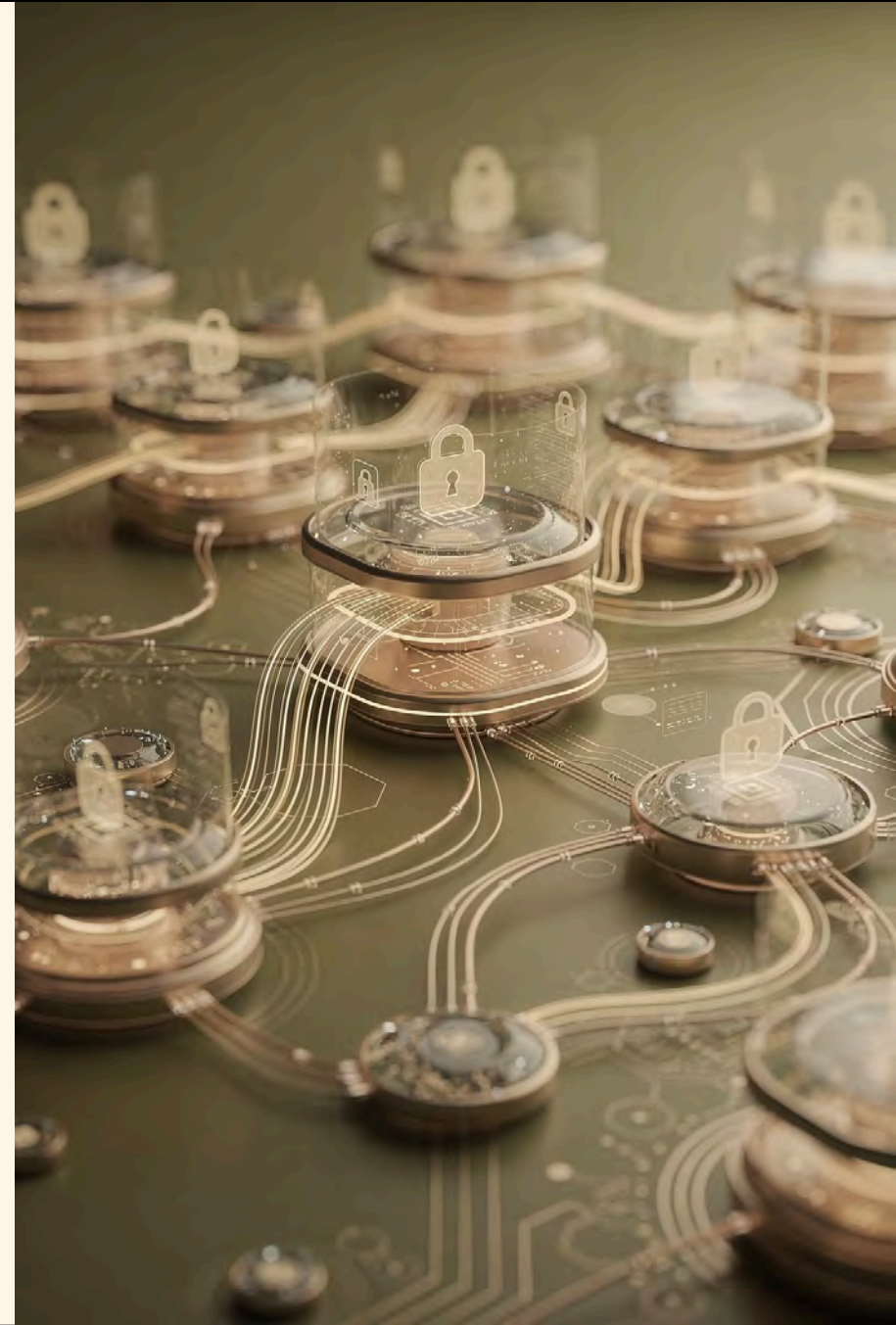
Measure and Iterate

Establish clear SLOs for availability, performance, and security then continuously improve against them



Building Robust, Secure, Scalable ML Systems

Production-grade ML platforms require intentional design across infrastructure, security, and organizational dimensions. By embedding reliability principles from the start and distributing responsibility across teams, organizations can build ML systems that deliver value continuously and securely at scale.



Key Takeaways

Foundational Resilience

Moving from monolithic legacy systems to cloud-native, distributed architectures on platforms like IBM Cloud SoftLayer.

High Availability

Achieved through multi-zone redundancy, automated failover, and stateless model serving for continuous ML operations.

Zero Trust Security

Segmentation protects ML workloads while maintaining development velocity and ensuring data integrity.

Performance Metrics

Track model serving uptime (99.9%+), training job completion rates (95%+), and incident response times for optimal performance.

Cross-functional Collaboration

Essential for production success, involving ML engineers, platform teams, and security specialists.

Hybrid Deployment

Balances cloud-native benefits with on-premises requirements for robust enterprise ML systems.


Thank You



Srinivas Talasila

SAP America Inc

Questions and discussion welcome. Let's connect on designing resilient ML infrastructure for your production systems.

 CONF42 ML