# Building Robust GPU Performance Validation into CI/CD Pipelines for AI Infrastructure

Presented By : Mohit Gupta

Intel Inc

Conf42 DevOps 2026.

# The Challenge: GPU Performance at Scale

Modern AI infrastructure demands systematic performance validation integrated directly into deployment workflows. Traditional post-deployment testing creates expensive feedback loops that slow innovation and increase costs.

Contemporary DevOps practices require shift-left approaches—detecting GPU performance bottlenecks during continuous integration before production deployment, not after.

# Why Performance Validation Matters

## Cost Impact

GPU resources represent significant infrastructure investment. Performance regressions directly translate to wasted compute cycles and increased operational expenses.

## Reliability Requirements

AI/ML workloads require predictable performance characteristics. Unexpected degradation impacts model training times, inference latency, and service-level agreements.

## Early Detection

Catching performance issues in CI/CD prevents production incidents, reduces mean time to resolution, and maintains deployment velocity for platform teams.

# Comprehensive GPU Validation Framework

Our framework embeds GPU microarchitecture validation into automated CI/CD workflows, specifically addressing compute-bound versus memory-bound workload characteristics critical for infrastructure reliability.

The approach establishes quantitative classification criteria, automated performance gates, and observable metrics that enable platform engineering teams to maintain high-performance standards throughout the deployment lifecycle.

# Understanding Workload Characteristics

### Compute-Bound Workloads

Characterized by high arithmetic logic unit (ALU) utilization with minimal memory system pressure:

- ALU utilization exceeding 70%
- Memory stall cycles below 20%
- Dominated by mathematical operations

### Memory-Bound Workloads

Limited by data transfer capabilities rather than computational throughput:

- Cache miss rates above 15%
- Memory bandwidth utilization exceeding 60%
- Sensitive to access patterns and coalescing

# Classification Criteria: Quantitative Metrics

## Compute-Bound Indicators

- **ALU Utilization:** Greater than 70% sustained usage

- **Memory Stall Cycles:** Less than 20% of execution time

- **Cache Hit Rates:** Typically 85% or higher for L1

## Memory-Bound Indicators

- **Cache Miss Rate:** Above 15% for working sets

- **Bandwidth Utilization:** Exceeding 60% of theoretical peak

- **Coalescing Efficiency:** Critical for structured vs. irregular access

# Systematic Benchmark Validation

## 01

### Dense Linear Algebra Operations

Matrix multiplication, factorization, and vector operations representing compute-intensive AI workloads.

## 02

### Irregular Memory Access Patterns

Graph traversal, sparse operations, and random access patterns common in recommendation systems and knowledge graphs.

## 03

### Hybrid Execution Scenarios

Mixed workload validation combining compute and memory operations to identify interaction effects.

# Performance Analysis: Compute-Bound Results

## 85%

### Sustained ALU Utilization

Achieved across representative dense linear algebra benchmarks

## 85%

### L1 Cache Hit Rate

Demonstrates effective data locality in compute-intensive operations

## <20%

### Memory Stall Cycles

Minimal memory system bottlenecks during computation phases

# Performance Analysis: Memory-Bound Results

## Bandwidth Utilization

Memory-bound workloads reached **75% bandwidth utilization**, demonstrating significant memory system pressure characteristic of data-intensive AI operations.

## Coalescing Efficiency Variance

- **Structured Access:** 90% coalescing efficiency for sequential patterns

- **Irregular Patterns:** 45% efficiency revealing optimization opportunities

# Critical Discovery: Mixed Workload Interactions

**Key Finding:** Concurrent execution of compute-bound and memory-bound workloads reveals performance degradation that isolated testing consistently misses.

## 1

### Isolated Testing

Individual workload characterization provides baseline performance expectations

## 2

### Interaction Effects

Resource contention and cache pressure during concurrent execution changes performance profiles

## 3

### Realistic Validation

Production environments require mixed workload testing to capture true performance characteristics

# Implementing Automated Validation as Code

### Define Performance Gates

Establish quantitative thresholds for ALU utilization, memory bandwidth, cache metrics, and coalescing efficiency within CI/CD configuration.

### Automate Benchmark Execution

Integrate representative workload suites into continuous integration workflows with consistent execution environments.

### Capture Observable Metrics

Collect GPU performance counters, memory access patterns, and utilization data for SRE observability and trend analysis.

# Platform Engineering Capabilities

## Prevent Performance Regressions

Automated gates block deployments that fail performance criteria, maintaining consistent service quality and preventing production incidents.

## Continuous Performance Visibility

Trend analysis and historical comparison enable proactive identification of performance drift and capacity planning.

## Shift-Left Optimization

Early detection reduces feedback loop time, accelerates development velocity, and minimizes cost of performance fixes.

# Building SRE Observability

## Key Observable Metrics

- GPU utilization trends across deployments

- Memory bandwidth saturation indicators

- Cache performance degradation alerts

- Workload classification distributions

These metrics enable SRE teams to establish service-level objectives, monitor infrastructure health, and respond proactively to performance anomalies.

# Key Takeaways for Platform Teams

**1** Embed Validation in CI/CD

Shift-left performance testing prevents costly production issues and accelerates deployment confidence for AI infrastructure.

**2** Classify Workload Characteristics

Quantitative metrics distinguish compute-bound from memory-bound workloads, enabling targeted optimization strategies.

**3** Test Mixed Workload Scenarios

Interaction effects during concurrent execution reveal performance characteristics that isolated testing cannot capture.

**4** Build Observable Systems

Performance metrics as code creates foundation for SRE practices, trend analysis, and continuous improvement.

# Thank You!

Mohit Gupta | Intel Inc

Conf42 DevOps 2026.