Beyond Uptime: Revolutionizing Fintech Reliability Through SRE Innovation

In today's landscape, organizations rapidly scaling their AI initiatives face critical challenges in efficiently managing Kubernetes resources. This presentation unveils battle-tested optimization strategies that directly impact both performance and cost efficiency, drawing from extensive experience with large-scale AI/ML deployments in production environments.

We'll explore advanced techniques that have consistently demonstrated dramatic improvements in resource utilization rates across diverse workloads, helping organizations achieve substantial cost reductions while maintaining high performance for critical ML operations.



Advanced GPU Optimization Techniques



Multi-Instance GPU (MIG) Partitioning

Physically partition NVIDIA A100/H100 GPUs to serve multiple workloads simultaneously, achieving up to 7x better GPU utilization in mixed inference workloads.

CUDA Memory Management

Implement custom memory allocation strategies to reduce fragmentation and optimize TensorCore usage, improving throughput by 15-30% for common vision models.

Precision Optimization

Systematically implement mixed and lower precision training/inference (FP16, INT8) with minimal accuracy loss, reducing memory footprint and increasing throughput by 2-4x.

4

These techniques have consistently demonstrated dramatic improvements in GPU utilization rates across diverse workloads, particularly for organizations running multiple models simultaneously on shared infrastructure.



Autoscaling Strategies for ML Workloads

Workload Profiling

000

{{i}}}

 \bigcirc

Analyze historical usage patterns across training and inference to establish baseline resource requirements and identify peak utilization periods.

Custom Metrics Pipeline

Implement Prometheus adapters to expose ML-specific metrics like queue depth, batch processing time, and GPU memory pressure to the Horizontal Pod Autoscaler.

Predictive Scaling

Deploy time-based scaling rules for recurring workloads and ML-driven predictive scaling for variable loads, reducing cold-start latency by up to 85%.

Buffer Management

Maintain optimal headroom in GPU clusters with deliberate over-provisioning during critical business hours, balanced against aggressive scaledown during low-demand periods.

Organizations implementing these sophisticated autoscaling strategies have achieved substantial cost reductions while maintaining high availability for critical inference workloads.

Multi-GPU Training Job Orchestration



The core of efficient ML operations centers on practical implementations that leverage both native Kubernetes capabilities and specialized tooling to efficiently handle multi-GPU training jobs. Organizations using these orchestration patterns have seen training job completion times improve by 30-50% through optimal resource placement and network topology awareness.

Advanced Monitoring and Observability



Organizations have successfully implemented advanced monitoring frameworks to achieve comprehensive observability across CPU, GPU, and memory utilization, enabling precise resource management and cost attribution. These systems allow ML engineers to quickly identify bottlenecks in training and inference pipelines.

The most successful implementations provide both real-time alerts and long-term trend analysis, supporting both immediate troubleshooting and strategic infrastructure planning.

Storage Optimization for ML Data Pipelines

Distributed File Systems

Implementation of specialized file systems like WekaFS and Lustre that have demonstrated up to 10x improvement in training data throughput compared to standard network storage.

- Parallel data access across hundreds of nodes
- Optimized for small file access patterns

Data Caching Layers

Deployment of in-memory caching services between persistent storage and compute to dramatically reduce I/O wait times for frequently accessed datasets.

- Reduced training startup times bv 65%
- Automatic eviction policies for optimal cache utilization

Storage Tiering

Implementing automated lifecycle policies to move data between hot, warm, and cold storage tiers based on access patterns and workload requirements.

- storage infrastructure
- unified namespace

These storage optimization techniques utilizing specialized file systems have demonstrated significant improvements in data access speeds, directly impacting the efficiency of ML training and inference operations.

40-60% cost reduction for Transparent access through



Cost Optimization with Spot Instances

Workload Classification

\$

 \otimes

E

θĨθ

Categorize ML workloads by fault tolerance, checkpoint requirements, and execution time to identify spot-compatible jobs

Fault Tolerance Implementation

Deploy automated checkpointing systems and stateful recovery mechanisms to gracefully handle preemption events

Bidding Strategy Optimization

Implement dynamic bidding strategies based on historical pricing data and current market conditions

Hybrid Deployment Models

Balance spot and on-demand instances with automated workload migration based on availability and pricing signals

Organizations leveraging these spot instance strategies have achieved cost savings of 60-80% compared to on-demand pricing, while maintaining acceptable performance for training workloads. The key to successful implementation is selecting appropriate workloads and implementing robust recovery mechanisms.

Optimal Node Pool Configurations



Heterogeneous Hardware Segmentation

Create specialized node pools for different GPU types (A100, V100, T4) with workload-specific taints and tolerations

Cost-Performance Balancing

Implement regular analysis to identify and adjust underutilized or overprovisioned node configurations

Optimized node pool configurations have helped organizations achieve remarkable cost savings while maintaining performance. The most successful implementations continuously evolve their node pool strategy based on workload telemetry and changing hardware options, rather than using a one-size-fits-all approach.

Resource Ratio Optimization Fine-tune CPU:GPU:Memory ratios to match specific ML workload profiles and prevent resource bottlenecks

Network Topology Alignment underlying network architecture for

Resource Quotas and Multi-Tenancy

Namespace Strategy

Implement logical separation between teams and workload types, using dedicated namespaces for production inference vs. experimental training.

- Role-based access controls per namespace
- Environment-specific configurations
- Isolated networking policies

Resource Controls

Deploy comprehensive quota management using ResourceQuotas and LimitRanges to prevent resource monopolization.

- GPU allocation limits per namespace
- Memory/CPU constraints with buffers
- Storage quota enforcement

Fair Scheduling

Implement priority classes and weighted fair queuing to ensure equitable access during resource contention periods.

Effective multi-tenancy has proven crucial for organizations supporting multiple ML teams on shared infrastructure. These approaches have helped organizations maintain equitable access to limited GPU resources while preventing "noisy neighbor" problems that can derail critical production workloads.

Production workload prioritization Guaranteed minimums for all teams Hierarchical resource distribution

Real-World Performance Gains



Cost Reduction

Average savings across surveyed organizations after implementing comprehensive optimization

2.8x

Training Throughput

Average improvement in training job completion time

67%

GPU Utilization

94%

Inference SLA

Increase in average cluster-wide GPU utilization

Average SLA achievement for production ML models

These metrics represent outcomes from organizations that successfully implemented the optimization strategies discussed throughout this presentation. Results vary based on initial infrastructure efficiency, workload characteristics, and implementation completeness.

The most dramatic improvements typically occur in environments that previously lacked ML-specific resource management approaches and were applying general-purpose Kubernetes patterns to specialized ML workloads.

Key Takeaways and Implementation Roadmap

Establish Baseline Metrics

Deploy comprehensive monitoring to understand current resource utilization patterns and identify the highest-impact optimization opportunities specific to your workloads.

- GPU/CPU utilization across workload types
- Training job completion times and variance
- Current infrastructure costs per model

Implement Quick Wins

Begin with optimization strategies that offer immediate returns with minimal disruption to existing workflows and infrastructure.

- Right-size resource requests and limits
- Enable basic node autoscaling
- Implement workloadappropriate storage classes

Deploy Advanced Optimizations

Gradually implement more sophisticated techniques to further enhance performance and efficiency.

- Custom metric-based autoscaling
- Spot instance integration
- Specialized node pools with topology awareness

These optimization techniques have consistently delivered substantial ROI across organizations of varying scales and workload complexities. The key to successful implementation is a methodical approach that balances immediate efficiency gains with long-term architectural improvements.

Continuous Refinement

Establish processes for ongoing optimization as workloads evolve and new technologies emerge.

Regular efficiency reviews

Automated cost anomaly detection

Infrastructure and model cooptimization

Thank you