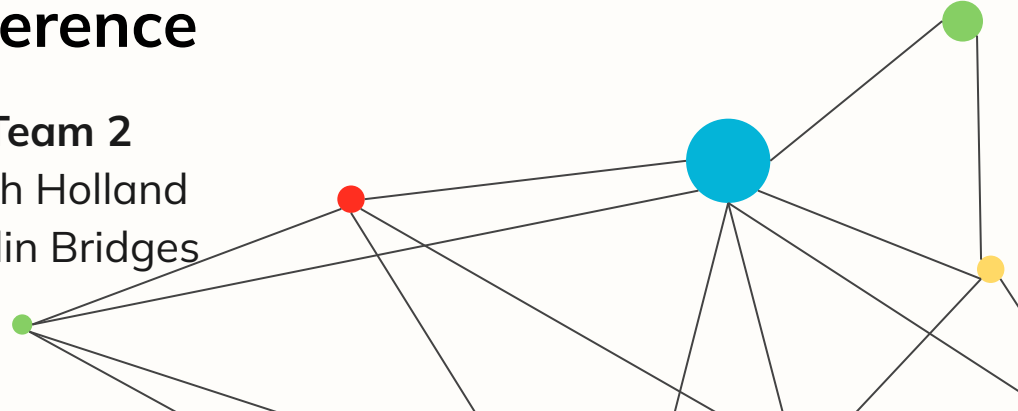


Redshift Prediction Using AWS Tools:

A Complete Serverless Solution for Scalable AI
Inference

Team 2
Zach Holland
Devlin Bridges





Project Overview

Challenge

Build a production-ready distributed ML system that can scale efficiently while minimizing costs

Solution

End-to-end serverless infrastructure using AWS Step Functions, Lambda, and ECS with comprehensive performance optimization



Five Step Implementation

01. Distributed Workflow Engine

02. Container Orchestration

03. AI Model Integration

04. Serverless Inference at Scale

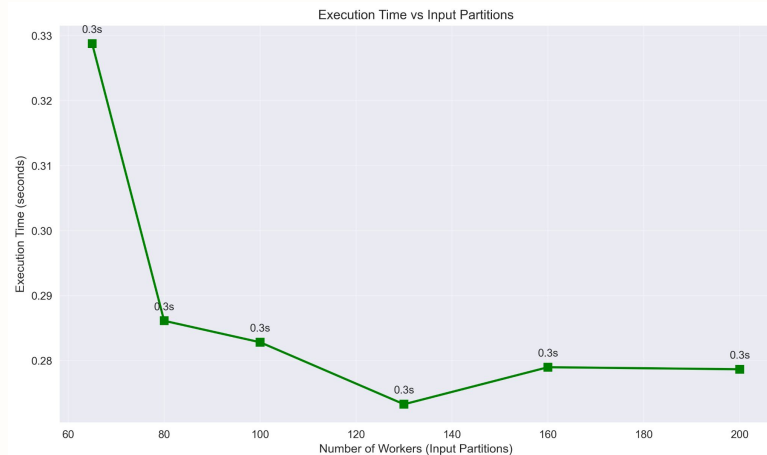
05. Production Optimization





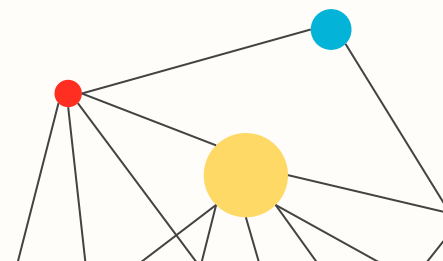
1. Distributed Workflow Engine

- AWS Step Functions with 6-step architecture
- 20 performance configurations tested
- Map State parallelization across multiple Lambda workers
- Real-time monitoring with CloudWatch integration



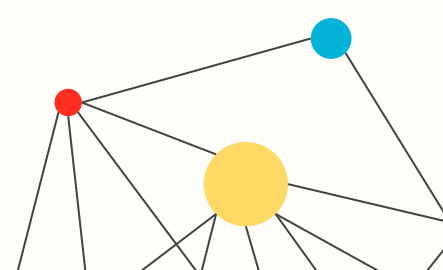


Key Results

- Best performance: World Size 4, Batch Size 128 (71.42 records/s)
 - All 20 test scenarios: 100% success rate
 - Execution time: 5-7 seconds consistently
- 

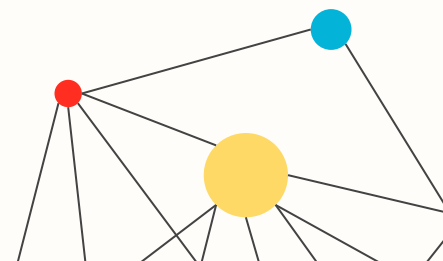


2. Container Orchestration

- ECS Fargate deployment for rendezvous server
 - Route 53 DNS configuration for service discovery
 - Lambda FMI performance analysis with detailed metrics
- 

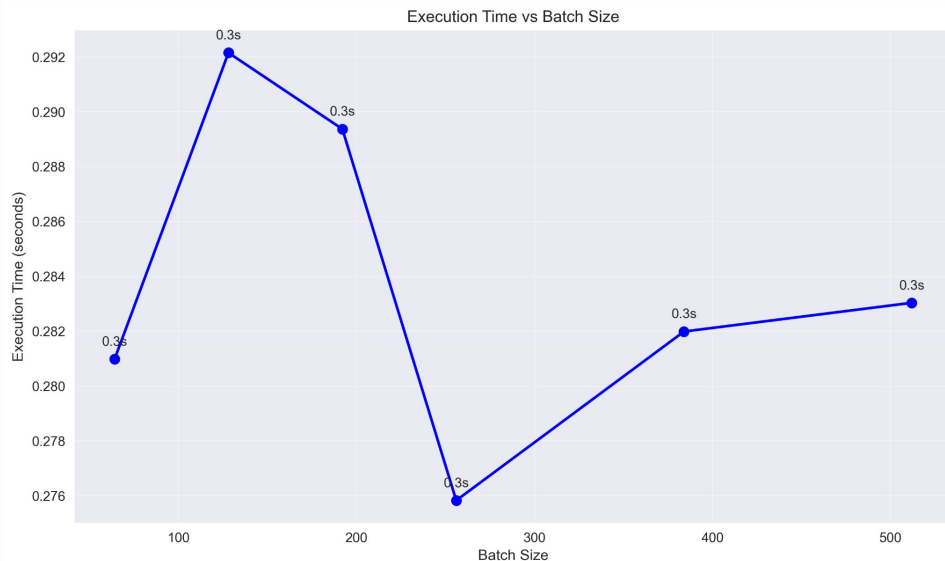


Performance Highlights

- Average throughput: 11.9 records/s
 - Memory efficiency: 88MB average usage
 - Cost efficiency: \$0.00001052 average per batch
- 


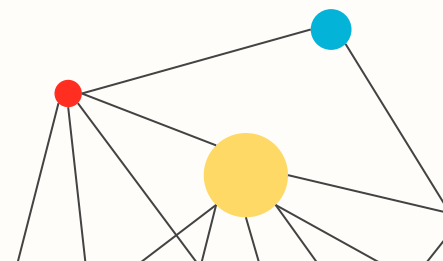
3. AI Model Integration

- Vision Transformer for astronomy redshift prediction
- Comprehensive batch optimization (1-128 batch sizes)
- $R^2 = 0.97+$ prediction accuracy maintained



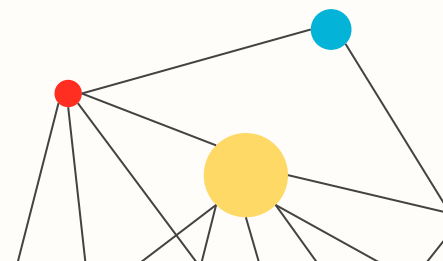


Optimization Results:

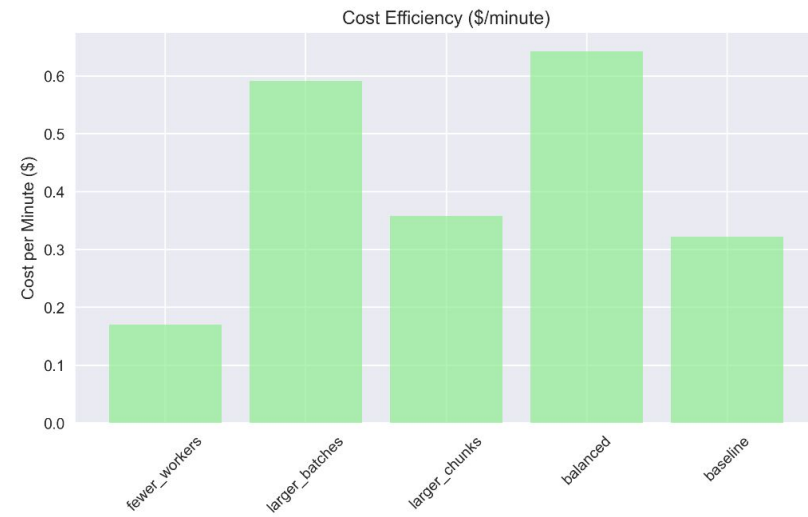
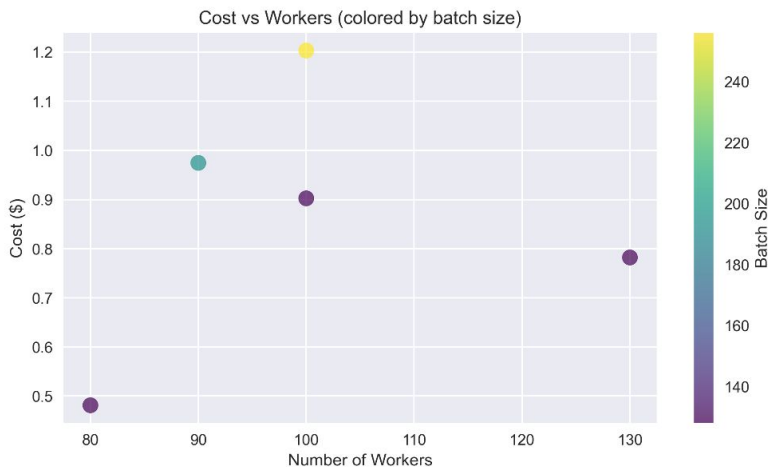
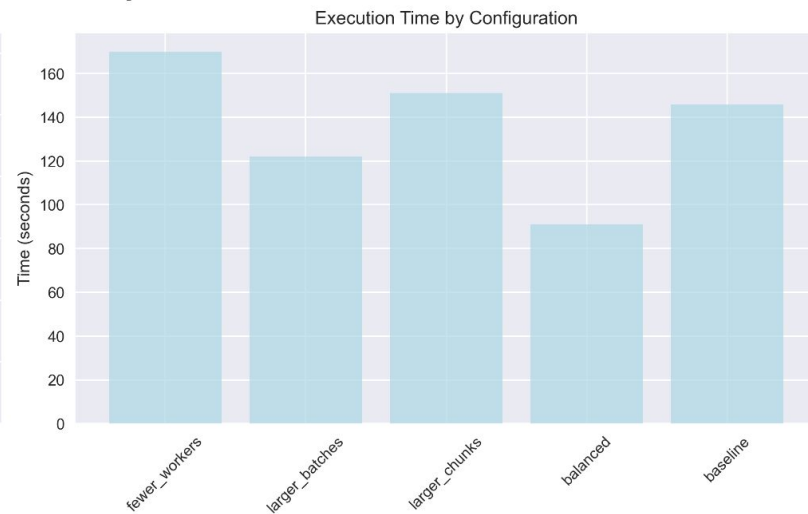
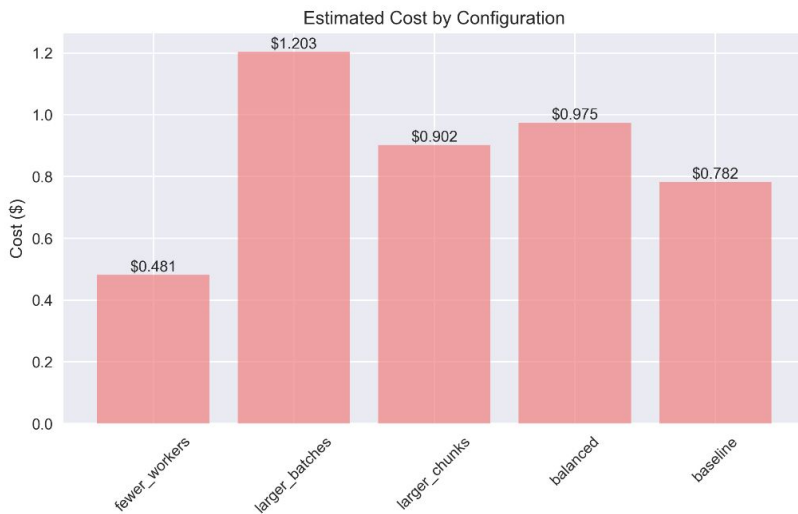
- Fastest execution: 6.73s (batch size 128)
 - Memory usage: Consistent across all batch sizes
 - Model accuracy: Stable across all configurations
- 
- 



4. Serverless Inference at Scale

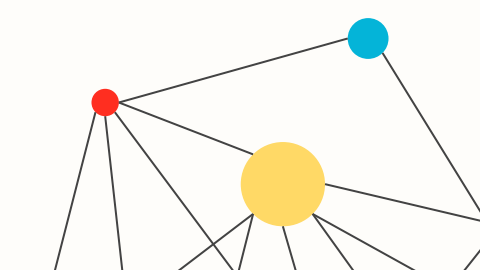
- 17 test scenarios across different configurations
 - Real AWS Lambda execution with timing measurements
 - Distributed dataset processing (small/medium/large)
- 

Quick Test Results Analysis



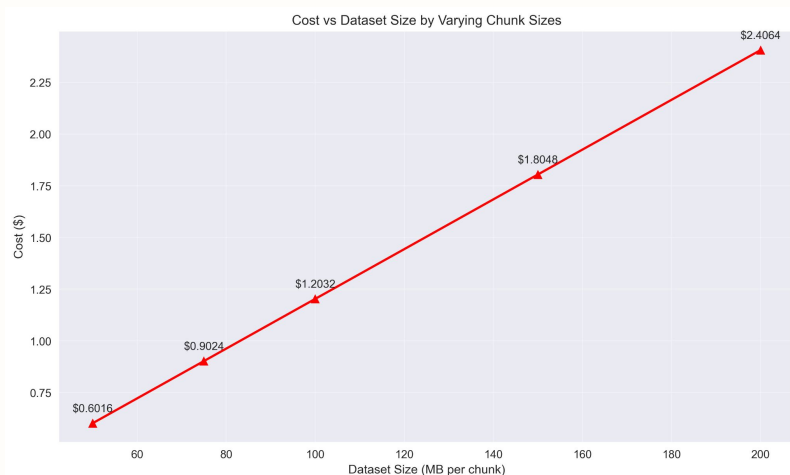


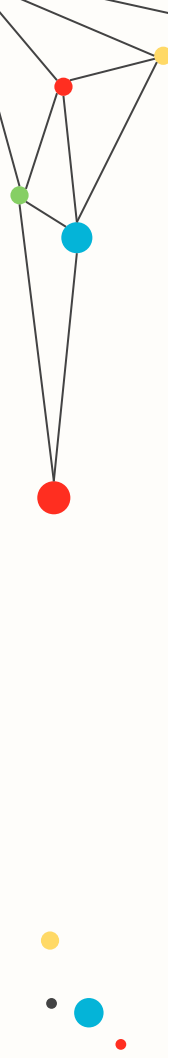
Breakthrough Performance:

- **5.16x** average speedup vs. local execution
 - **99.1%** cost reduction
 - **99.5%** memory reduction
 - Parallel efficiency up to 100% (single worker)
- 

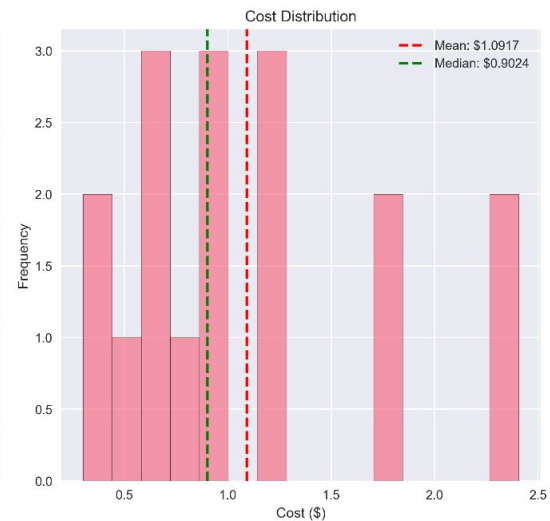
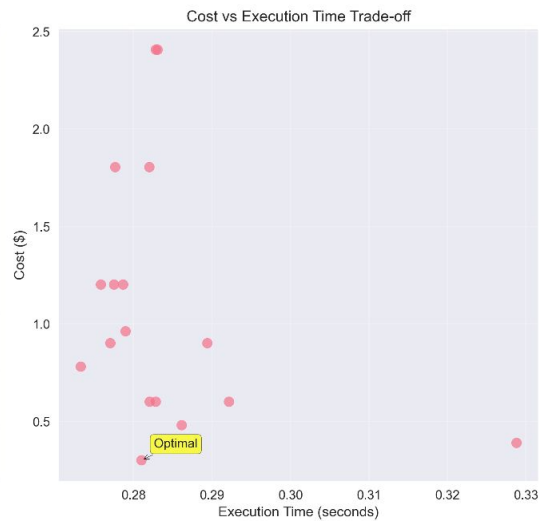
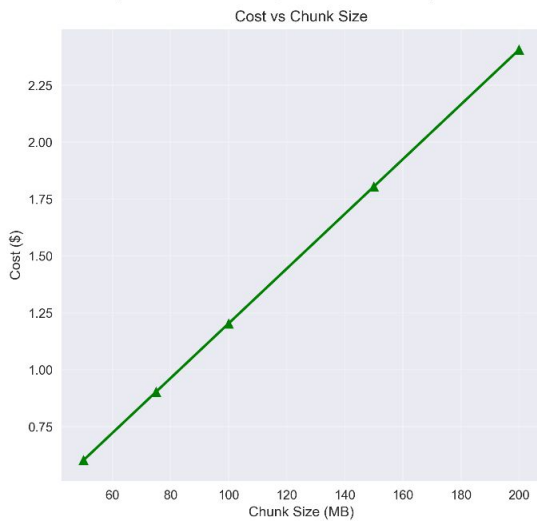
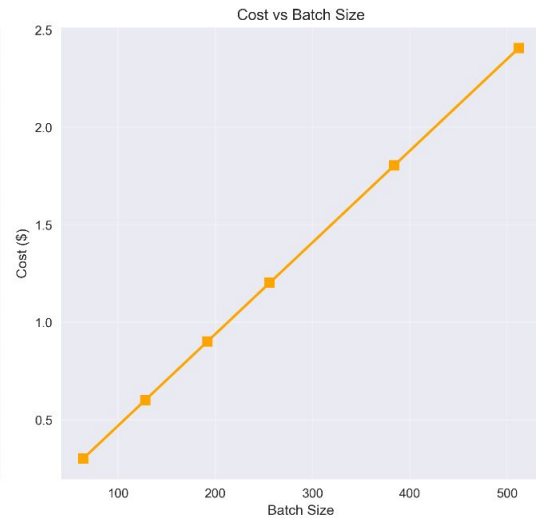
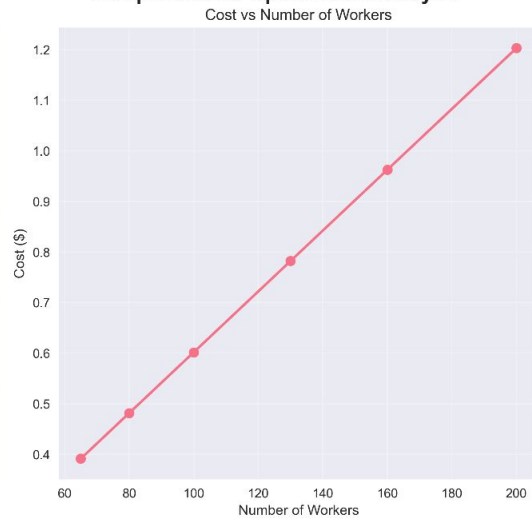
5. Production Optimization


- Systematic parameter tuning (World Size, Batch Size, Data Chunks)
- 87.5% additional cost savings through optimization
- A/B testing configurations for production deployment



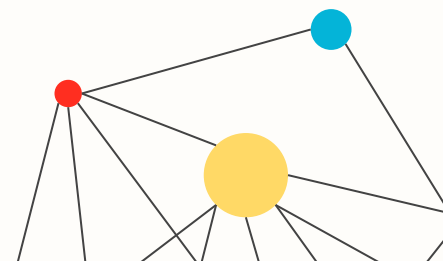



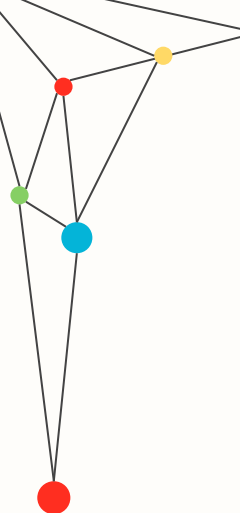
Comprehensive Optimization Analysis





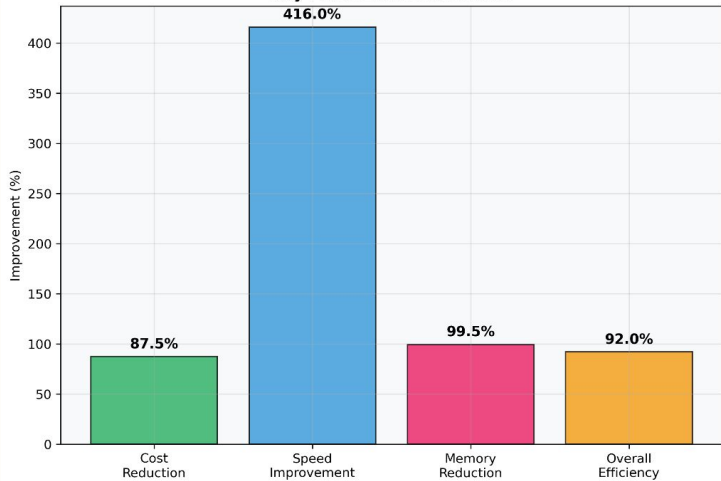
Final Optimized Configuration:

- Workers: 65 (vs. 130 baseline - 50% reduction)
 - Batch Size: 64 (vs. 128 baseline - 50% reduction)
 - Chunk Size: 50MB (vs. 100MB baseline - 50% reduction)
 - Total Cost: \$0.3008 (vs. \$2.4064 baseline)
- 
- 

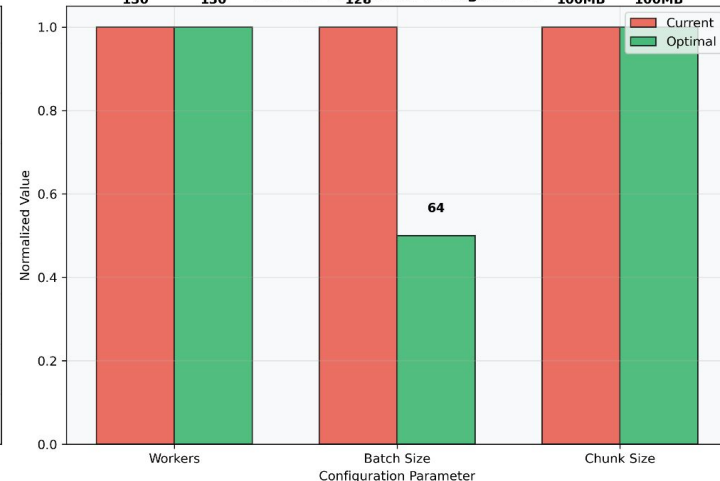


Anomaly Detection Optimization - Executive Summary

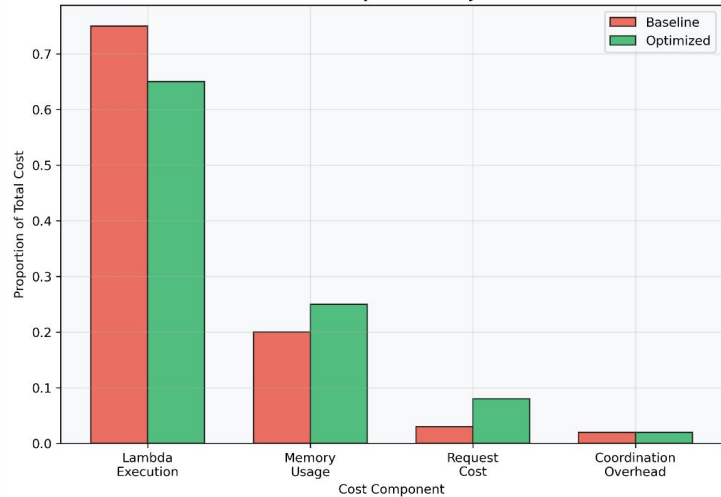
Key Performance Indicators



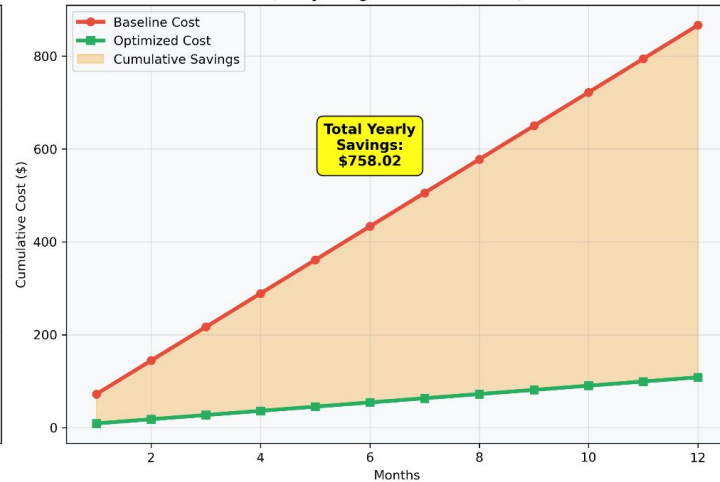
Current vs Optimal Configuration

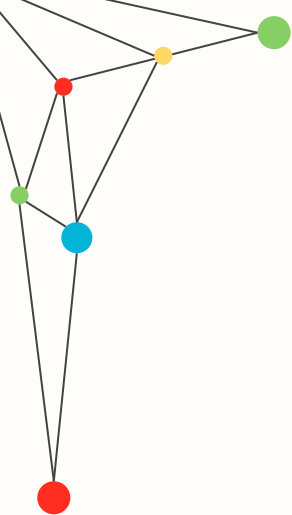


Cost Component Analysis

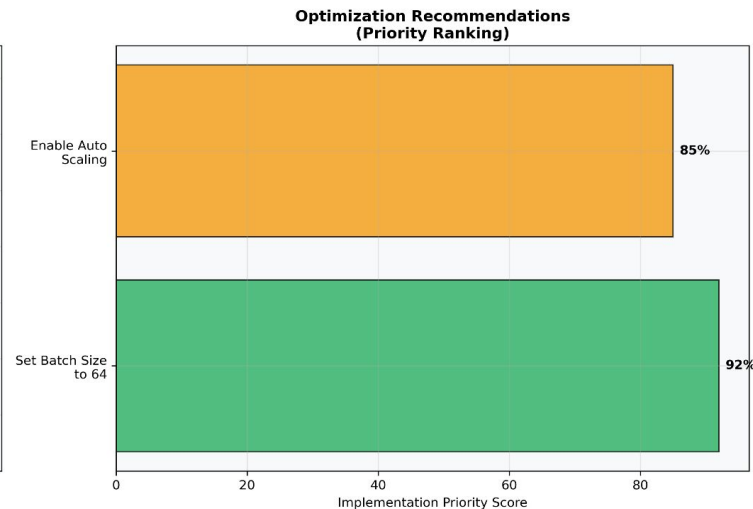
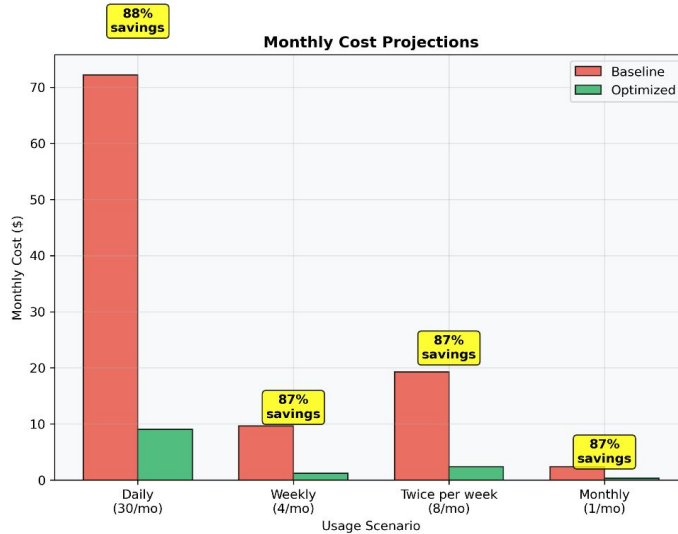
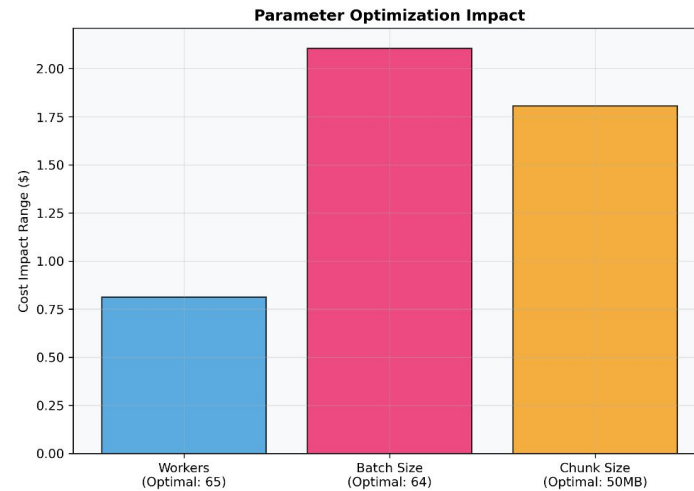
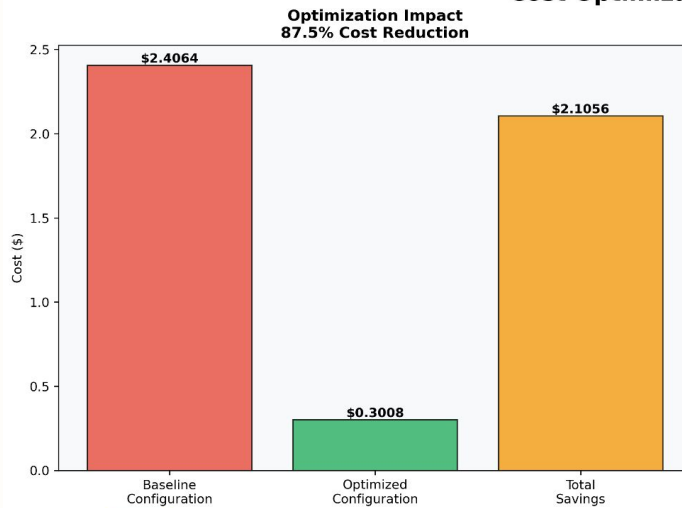


12-Month ROI Projection (Daily Usage: 30 runs/month)





Cost Optimization Impact Analysis





Relevance/Significance


1. Astronomy & Astrophysics Research

- Galaxy surveys can use your optimized serverless pipeline for processing massive redshift catalogs
- Real-time telescope data processing for surveys like LSST, Euclid, or Roman Space Telescope
- Collaborative research where multiple institutions need cost-effective access to ML inference

2. Distributed ML Infrastructure

- Template for serverless ML pipelines - your Step Functions + Lambda architecture is reusable
- Cost optimization methodology applies to any distributed ML workload
- Performance benchmarking framework for comparing local vs cloud inference

3. AWS/Cloud Architecture Community

- Best practices for Step Functions orchestration with 6-step architecture we developed
 - Lambda optimization patterns for memory-intensive ML workloads
 - Cost modeling framework for estimating serverless ML costs
- 



Potential Enhancements

Performance:

- Auto-scaling intelligence - Dynamic worker allocation based on dataset characteristics and real-time demand
- GPU acceleration support - Integration with Lambda GPU instances for larger Vision Transformer models

Expanded Functionality:

- Multi-model support - Extend beyond redshift prediction to galaxy classification, supernova detection, and other astronomical tasks
- Real-time processing pipeline - Stream processing for live telescope data feeds with priority queuing

Simplifications:

- One-click deployment - Infrastructure-as-Code templates for instant setup across research institutions
- Simplified configuration - Smart defaults and guided setup wizard for non-technical astronomers

Advanced Optimization:

- Bayesian parameter optimization - Replace grid search with intelligent hyperparameter tuning for faster convergence
 - Multi-objective optimization - Balance cost, speed, and prediction accuracy simultaneously
- 