# Cloud Computing Project Report: Autoscaler Performance Comparison

## Group 17 Members

1. Talal Ahmed
2. Zain Khalid
3. Usama Arif
4. Muhammad Abdullah Bin Saif

## Objective

The goal of this project was to evaluate a custom autoscaler against Kubernetes Horizontal Pod Autoscaler (HPA) using two CPU utilization targets: 70% and 90%. The aim was to reduce latency and CPU usage while maintaining performance.

## Experiment Setup

We ran three experiments:

- HPA 70% Target – Kubernetes scales pods when average CPU exceeds 70%.
- HPA 90% Target – Scaling happens when CPU usage goes beyond 90%.
- Custom Autoscaler – Our tailored solution that adjusts based on optimized thresholds and workload patterns.

For each setup, we tracked:

- 99th percentile latency: to measure worst-case performance.
- CPU usage: to understand how efficiently resources were used.

## Results Summary

| Metric | HPA 70% | HPA 90% | Custom Autoscaler |
|---|---|---|---|
| Final 99th % latency (sec) | ~ 0.33 | ~ 0.39 | ~ 0.31 |
| Final CPU usage (%) | ~ 6.06% | ~ 3.03% | ~ 6–11% dynamic |

# Observations

1. **Latency Comparison**
   - The custom autoscaler consistently maintained lower latency, especially under high load.
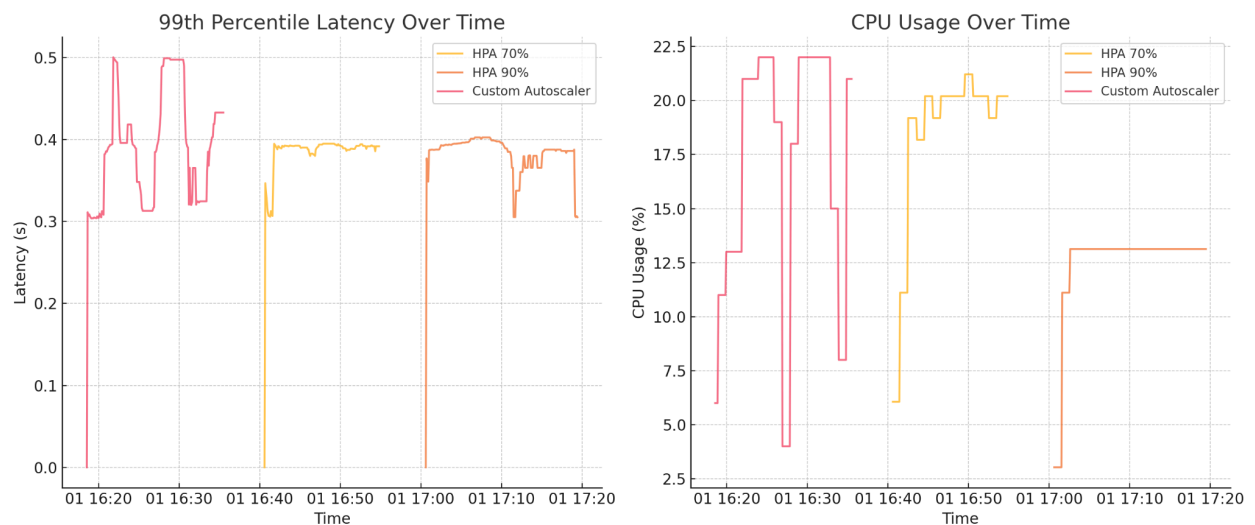   - HPA 90% suffered from occasional latency spikes due to delayed scaling decisions.
2. **CPU Usage**
   - While HPA 90% used slightly less CPU, it came at the cost of higher latency.
   - The custom autoscaler dynamically adjusted usage, striking a better balance between performance and efficiency.
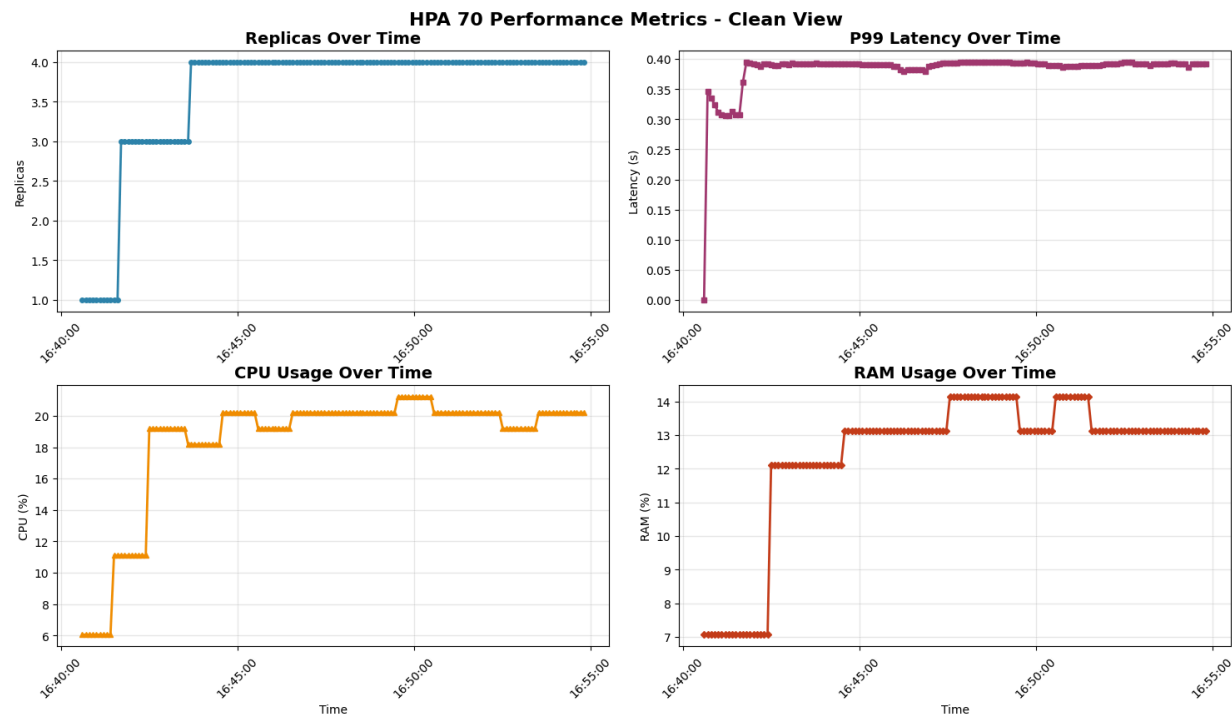3. **Overall Performance**
   - The custom autoscaler outperformed both HPA configurations in terms of latency.
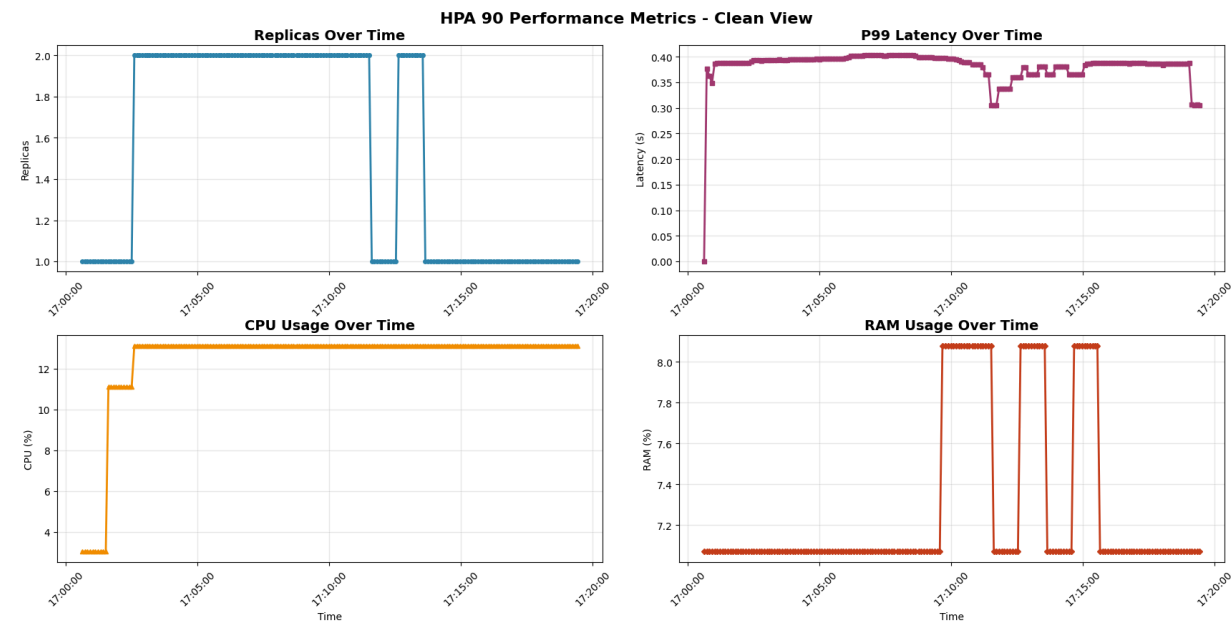   - It demonstrated intelligent scaling with moderate CPU usage, adapting better to real-time load patterns.

# Graphics

# HPA 70

## HPA 70 Performance Metrics - Clean View

### Replicas Over Time



### P99 Latency Over Time



### CPU Usage Over Time



### RAM Usage Over Time



# HPA 90

## HPA 90 Performance Metrics - Clean View

### Replicas Over Time



### P99 Latency Over Time



### CPU Usage Over Time



### RAM Usage Over Time

# Custom Autoscaler

## Custom Autoscaler Performance Metrics - Clean View



### Replicas Over Time

### P99 Latency Over Time

### CPU Usage Over Time

### RAM Usage Over Time