



# 🚀 KUBERNETES HPA MASTERY

Scale Your Apps Like a Pro

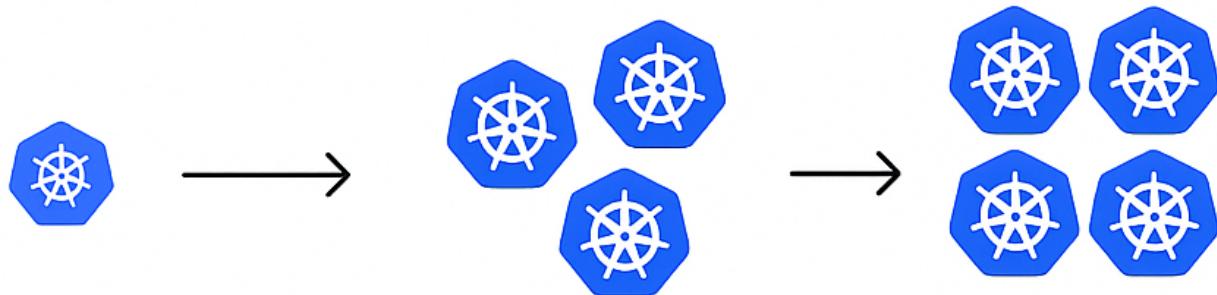
- ✓ Auto-scaling fundamentals
- ✓ Real-world examples
- ✓ Best practices & troubleshooting
- ✓ Production-ready configs

Perfect for DevOps &  
Platform Engineers

👉 Swipe to learn more



# What is Horizontal Pod Autoscaler?



HPA automatically scales your Kubernetes pods based on:

- CPU utilization
- Memory usage
- Custom metrics

HPA = Smart resource management



## Key Benefits:

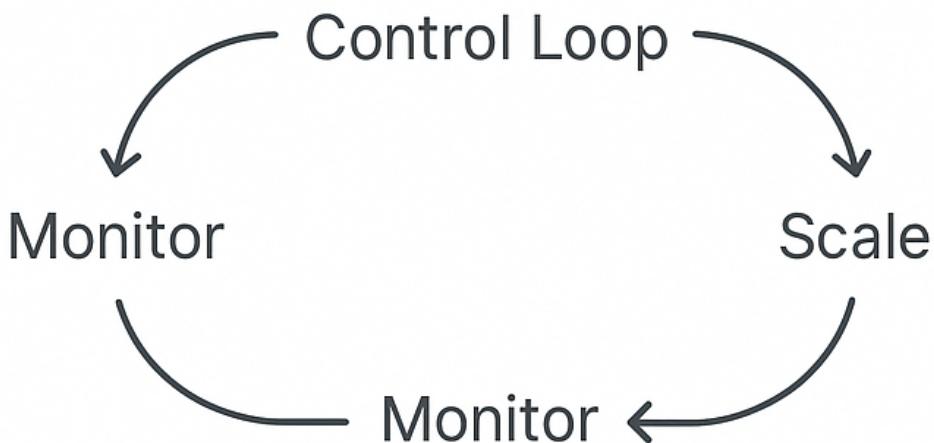
- ✓ Zero manual intervention
- ✓ Cost optimization
- ✓ Performance consistency
- ✓ Handle traffic spikes automatically

**HPA = Smart resource management**

# How HPA Works Under the Hood



- ① Metrics Server collects data
- ② HPA evaluates thresholds
- ③ Adjusts replica count in Deployment
- ④ Deployment Controller creates/deletes pods



-  HPA doesn't create pods directly!  
It modifies the Deployment's replica count

# HPA Prerequisites Checklist

-  Metrics Server running in cluster
-  CPU/Memory requests defined in pods
-  Target workload (Deployment/StatefulSet)
-  Proper RBAC permissions
-  **Common Mistake:**  
Missing resource requests =  
HPA won't work!
-  **Example:**  
resources:  
  requests:  
    cpu: 100m  
    memory:128Mi



## Basic HPA Configuration

```
apiVersion: autoscaling/v2
kind: HorizontalPodAutoscaler
metadata:
  name: nginx-hpa
spec:
  scaleTargetRef:
    apiVersion: apps/v1
    kind: Deployment
    name: nginx-deployment
  minReplicas: 2
  maxReplicas: 10
metrics:
  type: Resource
  resource: cpu
  target: Utililzation
  averageUtilization: 50
```



This scales between 2-10 pods at 50% CPU

# Quick HPA Setup with kubectl

```
# Create HPA in one command  
kubectl autoscale deployment nginx \  
  --cpu-percent=50 \  
  --min=1  
  --max=5  
  
# Check HPA status  
kubectl get hpa  
  
# Watch real-time scaling  
kubectl get hpa -w
```

 Pro Tip: Great for testing!  
For production, use YAML manifests

# Beyond CPU: Advanced Metrics



## Resource Metrics

- CPU utilization
- Memory utilization



## Custom Metrics

- Queue depth
- Request latency
- Active connections
- Business metrics



## External Metrics

- CloudWatch metrics
  - Prometheus metrics
  - Third-party APIs
- CloudWatch metrics
  - Prometheus metrics
  - Third-party APIs



## External Metrics



Use KEDA for complex custom metrics!



# Multi-Metric HPA Example

metrics:

- type:Resource  
resource:cpu  
target:  
    type:Utilization  
    averageUtilization:70
- type:Resource  
resource:memory  
target:  
    type:Utilization: 80



## HPA Logic:

Scales based on the metric requiring the HIGHEST number of replicas

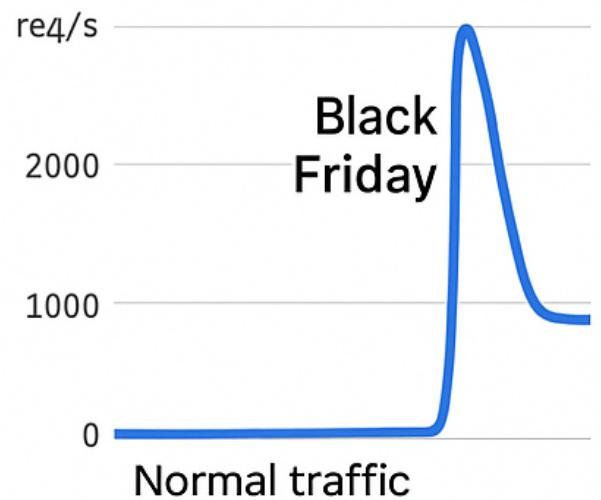
CPU needs 3 pods, Memory needs 5 pods

→ HPA scales to 5 pods

# 🛒 Real-World Example: E-commerce Black Friday

## 📈 Scenario:

- Normal traffic: 100 req/s (2 pods) 🛒
- Black Friday: 2000 req/s spike 🛒



## ⚙️ HPA Configuration:

- minReplicas: 2
- maxReplicas: 50
- targetCPU: 60%

## 📈 Results:

- Autoscaled to 32 pods in 10 minutes
- Handled 20x traffic increase
- Zero downtime
- 40% cost savings vs fixed scaling

⌚ Key: Proper resource requests + testing



# HPA Best Practices



## DO:

- Set resource requests accurately
- Test scaling behavior in staging
- Monitor scaling events
- Use appropriate min/max replicas
- Set up proper alerting



## DON'T:

- Skip resource requests
- Set aggressive thresholds
- Ignore scaling delays
- Run HPA + VPA together
- Forget about node capacity



Start conservative,  
optimize gradually



# Common HPA Issues & Fixes



**Problem:** HPA shows 'Unknown' metrics

**Solution:** Check Metrics Server is running



**Problem:** No scaling happening

**Solution:** Verify resource requests set



**Problem:** Constant scaling up/down

**Solution:** Adjust thresholds & delays



**Problem:** Pods pending during scale-up

**Solution:** Check node capacity/ Cluster Autoscaler



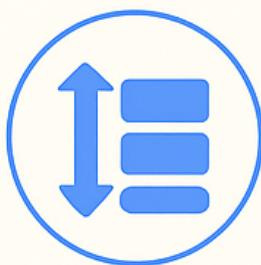
**Problem:** Slow scaling response

**Solution:** Review controller delays



**Debug command:** kubectl describe hpa

# ⚖️ Autoscaling Types Comparison

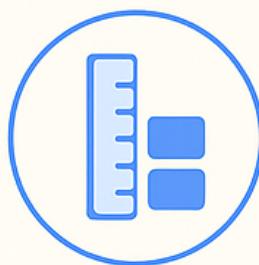


## HPA (Horizontal):

Scales  
NUMBER of pods

Best for:  
Stateless apps

Handles:  
Traffic spikes

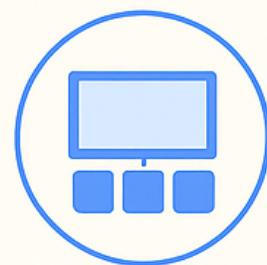


## VPA (Vertical):

Scales  
POD RESOURCES

Best for:  
Right-sizing

Handles:  
Resource optimization



## CAS (Cluster):

Scales  
NUMBER of nodes

Best for:  
Node management

Handles:  
Cluster capacity

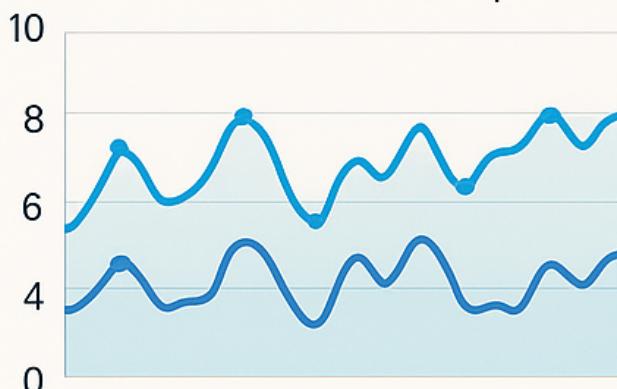


They work together for complete elasticity!

# Production Monitoring

## Monitor Your HPA in Production

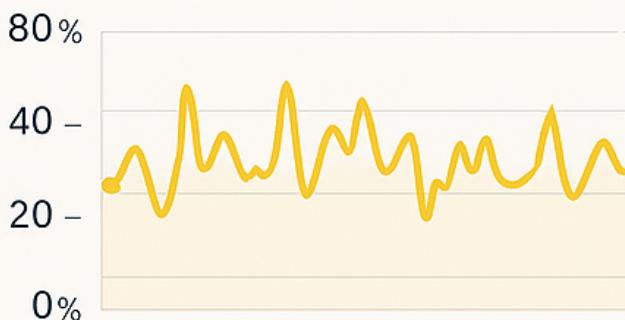
Current vs desired replicas



### Key Metrics to Track:

- Current vs desired replicas
- Scaling frequency
- Resource utilization trends
- Scale-up/down latency

CPU Utilization



### Essential Alerts:

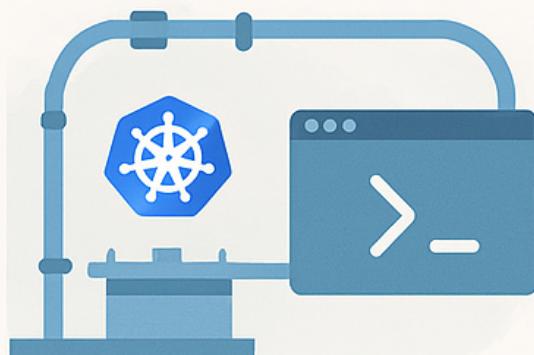
- HPA unable to scale
- Frequent scaling events
- Max replicas reached
- Metrics server unavailable

4.2s



### Tools:

- Prometheus + Grafana
- Kubernetes Dashboard
- kubectl top pods



Scaling events = kubectl get events



# Pro HPA Configuration Tips



## Scaling Delays

```
--horizontal-pod-autoscaler--  
downscale-delay=5m  
--horizontal-pod-autoscaler--  
upscale-delay=3m
```



## Stabilization Windows

```
behavior:  
  scaleUp:  
    stabilizationWindowSeconds: 60  
  scaleDown:  
    stabilizationWindowSeconds: 390
```



## Advanced Targeting

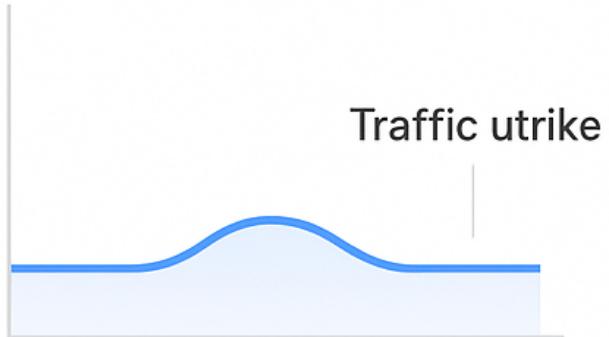
```
# Scale based on requests per second  
type: Pods  
metric: name:requests_per_second  
target: AverageValue
```

# 🚀 HPA Performance Optimization

## Before HPA



## After HPA



- Manual scaling decisions
- Over-provisioned resources
- Poor cost efficiency
- Traffic spike failures

- 60% cost reduction
- 99.9% uptime during spikes
- Automatic resource optimization
- Zero manual intervention

## ⌚ Real Results:

- Netflix: Handles millions of requests
- Spotify: 40 % resource savings
- Airbnb: 50 % faster response times

💰 ROI: HPA pays for itself in months!



# Ready to Scale Like a Pro?

## What's Your Next Step?



### Start Today:

1. Implement basic CPU-based HPA
2. Add resource requests to your pods
3. Monitor and optimize thresholds
4. Expand to custom metrics



Questions? Drop them in comments!



Share if this helped you



Tag a DevOps friend who needs this



Want more Kubernetes content?

Follow me for weekly platform engineering tips!

#Kubernetes #DevOps #CloudNative  
#Scaling #PlatformEngineering #SRE  
#CloudComputing