

AWS Kubernetes Lab Guide: HPA, Cluster Autoscaler, Prometheus & Grafana

AWS Kubernetes Class Notes: Horizontal Pod Autoscaling (HPA)

Recap: Previous Class

- Deployed an application (e.g., Node.js app) on Kubernetes.
 - Connected the app to a domain purchased from GoDaddy.
-  Request Flow:

User → Load Balancer → Ingress → Service → Deployment → Pods

Clone Kubernetes Autoscaling Manifests

GitHub: <https://github.com/sibasish934/kubernetes-manifests.git>

Navigate to: auto-scaling

What is HPA (Horizontal Pod Autoscaler)?

Automatically scales pods based on observed CPU, memory, or custom metrics.

Step-by-Step Lab Implementation

- 1 Cluster Creation: EKS cluster already created with nodes and pods running.
- 2 Load Generation: Use BusyBox pod to simulate CPU load.
- 3 Check Resource Usage: Use `kubectl top pods`.
- 4 Configure HPA: Define CPU threshold in hpa.yaml.

- 5** Apply HPA: `kubectl apply -f hpa.yaml`
- 6** Monitor Scaling: Use `kubectl get hpa` and `kubectl get pods --watch`
- 7** Trigger Load Again: Observe scaling up and down.

Conclusion

- HPA optimizes performance by auto-scaling pods.
- Requires metrics-server.
- Can be extended to memory or custom metrics.

Cluster Autoscaler Setup (EKS)

-  Attach IAM Policy to user
-  Set up OIDC provider using eksctl
-  Create IAM role for autoscaler
-  Update YAML with role ARN & cluster name
-  Apply YAML
-  Check logs with `kubectl logs`
-  Verify in AWS Console (Auto Scaling Groups)
-  Optional: Apply additional test deployments.

Clean Up

Delete HPA test deployment and service:

- `kubectl delete -f hpa-testing-deploy.yaml`
- `kubectl delete -f hpa-testing-service.yaml`

Monitoring with Prometheus & Grafana

-  Prometheus: collects metrics
-  Grafana: displays dashboards

Monitoring Setup Lab

- 1** Check Metrics API: `kubectl get --raw /metrics`
- 2** Ensure AWS Load Balancer Controller is installed
- 3** Install EBS CSI Driver: for persistent volumes

Installing Prometheus

- Create `prometheus` namespace
- Add Helm repo: prometheus-community
- Install with Helm and gp2 storage class
- Verify pods
- Port-forward to view Prometheus UI at <http://localhost:9090>

Installing Grafana

- Create `grafana` namespace
- Add Helm repo: grafana
- Install via Helm (set admin password)
- Port-forward to <http://localhost:3000>
- Login: admin/admin123
- Import Dashboard ID 1860 for node metrics


```

PS C:\users\Ayush Singh\downloads\kubernetes-manifests\auto-scaling> kubectl apply -f cluster-auto-scaler.yaml
serviceaccount/cluster-autoscaler created
clusterrole.rbac.authorization.k8s.io/cluster-autoscaler created
role.rbac.authorization.k8s.io/cluster-autoscaler created
clusterrolebinding.rbac.authorization.k8s.io/cluster-autoscaler created
rolebinding.rbac.authorization.k8s.io/cluster-autoscaler created
deployment.apps/cluster-autoscaler created
PS C:\users\Ayush Singh\downloads\kubernetes-manifests\auto-scaling> kubectl get pods -n kube-system
NAME                               READY   STATUS    RESTARTS   AGE
aws-node-6nk8f                     2/2     Running   0          51m
aws-node-n97mx                     2/2     Running   0          51m
cluster-autoscaler-64b685b46b-75529 1/1     Running   0          27s
coredns-6799d65cb-fncjt            1/1     Running   0          53m
coredns-6799d65cb-v4s76            1/1     Running   0          53m
kube-proxy-cnkcj                   1/1     Running   0          51m
kube-proxy-wptf4                   1/1     Running   0          51m
metrics-server-6c7659898f-4679w    1/1     Running   0          53m
metrics-server-6c7659898f-wwq29    1/1     Running   0          53m



⌚ Auto Scaling group updated successfully
×



Auto Scaling groups (1/1) Info
Last updated 1 minute ago 
Launch configurations
Launch templates
Actions ▾
Create Auto Scaling group



| <input checked="" type="checkbox"/>   Name                                                                              | Launch template/configuration        | Instances | Status | Desired capacity | Min | Max |
|-------------------------------------------------------------------------------------------------------------------------|--------------------------------------|-----------|--------|------------------|-----|-----|
| <input checked="" type="checkbox"/> eks-standardworkers-e2cbe053-53a9-6cc7-fb72-2a389c25<br>53a9-6cc7-fb72-2a389c257ee4 | eks-e2cbe053-53a9-6cc7-fb72-2a389c25 | 2         | -      | 2                | 1   | 5   |



```

PS C:\users\Ayush Singh\downloads\kubernetes-manifests\auto-scaling> helm install aws-load-balancer-controller eks/aws-load-balancer-controller -n kube-system --set clusterName=test-cluster --set serviceAccount.create=false --set serviceAccount.name=aws-load-balancer-controller --set region=ap-south-1 --set vpcId=@646886ff346a9488
NAME: aws-load-balancer-controller
LAST DEPLOYED: Mon Jun 30 18:27:19 2025
NAMESPACE: kube-system
STATUS: deployed
REVISION: 1
TEST SUITE: None
NOTES:
AWS Load Balancer controller installed!
PS C:\users\Ayush Singh\downloads\kubernetes-manifests\auto-scaling> eksctl create addon --name aws-ebs-csi-driver --cluster dev-cluster --region ap-south-1 --service-account-role-arn arn:aws:iam::26985572870:role/eks-cluster-autoscaler-role --force
2025-06-30 18:29:49 [i] Kubernetes version "1.32" in use by cluster "dev-cluster"
2025-06-30 18:29:49 [i] IRSA is set for "aws-ebs-csi-driver" addon; will use this to configure IAM permissions
2025-06-30 18:29:49 [i] the recommended way to provide IAM permissions for "aws-ebs-csi-driver" addon is via pod identity associations; after addon creation is completed, run `eksctl utils migrate-to-pod-identity`
2025-06-30 18:29:49 [i] using provided ServiceAccountRoleARN "arn:aws:iam::26985572870:role/eks-cluster-autoscaler-role"
2025-06-30 18:29:49 [i] creating addon: aws-ebs-csi-driver

Welcome to Grafana

✉

🔒

Forgot your password?

Documentation | Support | Community | Open Source | v8.3.0 (30bb7a93ca)


```


```