## **Kubernetes Autoscaler**



# By Ritesh Kumar Singh

Email Address: - riteshkumarsingh9559@gmail.com

LinkedIn: - https://www.linkedin.com/in/ritesh-kumar-singh-41113128b/

GitHub: - https://github.com/singhritesh85



या कुन्देन्दुतुषारहारधवला या शुभ्रवस्त्रावृता या वीणावरदण्डमण्डितकरा या श्वेतपद्मासना। या ब्रह्माच्युत शंकरप्रभृतिभिर्देवैः सदा वन्दिता सा मां पातु सरस्वती भगवती निःशेषजाड्यापहा॥

## 1. First Section

#### **Kubernetes Autoscaler**

#### 1. EKS Cluster Autoscaler

As the name indicates EKS cluster autoscaler will scale the EKS Node-Group nodes between the minimum and maximum number of nodes of EKS Cluster. For the current scenario I will allow the nodes of EKS Node-Group to scale between 1 to 4. You can change it as per your project requirement.

```
scaling_config {
  desired_size = 1
  max_size = 4
  min_size = 1
}
```

Always remember the desired size can not be less than minimum size of the Autoscaling Group.

I used the terraform script present in the GitHub Repo

https://github.com/singhritesh85/Autoscale.git at the path Autoscale/Autoscale-EKS/Autoscale-EKS-using-cluster-autoscaler to create the EKS Cluster and after creation of EKS Cluster run the manifests file cluster-autoscaler.yaml as shown in the screenshot attached below.

```
terraform-eks-withaddons-cluster-autoscaler]# kubectl apply -f cluster-autoscaler.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
                terraform-eks-withaddons-cluster-autoscaler]# kubectl get pods -n kube-system
NAME
                                  READY
                                         STATUS
                                                  RESTARTS
                                                           AGE
                                         Running
cluster-autoscaler-
                                         Running
                                                            53s
                                                            118s
                                         Running
coredns-
                                                            118s
ebs-csi-controller
                                         Running
ebs-csi-controller
                                         Running
                                                            1145
ebs-csi-node-
                                                            114s
                                         Running
                                         Running
kube-proxy-
                                                            2m43s
                                         Running
                                         Running
                                        cluster-autoscaler]# kubectl logs -f cluster-autoscaler-
```

As per the requirement new nodes will be added and before the Scale down it will wait for 10 minutes and then scale down.

When you will check the logs of cluster-autoscaler pod present in kube-system namespace then you can check how the scale up and scale down will be going on.

For the demonstration purpose I created a deployment named as dexter with number of pods required is 20 and I checked some of the pods are in running condition and other are in pending condition due to unavailability of compute resource.

```
[root@ ~]# kubectl apply -f deployment.yaml
deployment.apps/dexter created
```

```
[root@ ~]# cat deployment.yaml
apiVersion: apps/v1
kind: Deployment
metadata:
 name: dexter
spec:
 selector:
   matchLabels:
     app: nginx
 replicas: 20
 template:
   metadata:
     labels:
       app: nginx
   spec:
     containers:
      - name: demo
       image: nginx
       ports:
        - containerPort: 80
```

[root@	~]#	kubectl	get pods		
NAME		READY	STATUS	RESTARTS	AGE
dexter-		0/1	Pending	0	33s
dexter-		1/1	Running	0	33s
dexter-		0/1	Pending	0	33s
dexter-		0/1	Pending	0	33s
dexter-		1/1	Running	0	33s
dexter-		0/1	Pending	0	33s
dexter-		0/1	Pending	0	33s
dexter-		0/1	Pending	0	33s
dexter-		0/1	Pending	0	33s
dexter-		1/1	Running	0	33s
dexter-		0/1	Pending	0	33s
dexter-		0/1	Pending	0	33s
dexter-		0/1	Pending	0	33s
dexter-		0/1	Pending	0	33s
dexter-		0/1	Pending	0	33s
dexter-		1/1	Running	0	33s
dexter-		0/1	Pending	0	33s
dexter-		0/1	Pending	0	33s
dexter-		1/1	Running	0	33s
dexter-		1/1	Running	0	33s

Then I checked the logs of cluster-autoscaler pod present in the kube-system namespace and found that two new nodes had been created and pods are in pending state launched on the nodes of the EKS Cluster.

[root@	~]# kubectl	get nodes								
NAME			STATUS	ROLES	AGE	VERSIO	N			
ip-	.us-east-2.compute.	internal	Ready	<none></none>	7m11s	v1.30.	4-eks-a737599			
ip-	us-east-2.compute.i	nternal	Ready	<none></none>	17m	v1.30.	4-eks-a737599			
ip-	.us-east-2.compute.	internal	Ready	<none></none>	7m14s	v1.30.	4-eks-a737599			
[root@	~]# kubectl	get pods	o wide							
NAME	READY	STATUS	RESTARTS	AGE	IP		NODE		NOMINATED NODE	READINESS GATES
dexter-	1/1	Running	0	8m22s	10.	131		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	245		us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	246		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	51		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	20		us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	238		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	190		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	215		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	136		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	91		us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	79		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	97		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	18		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	45		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	.98		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	52		us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	189		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	139		.us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	53		us-east-2.compute.internal	<none></none>	<none></none>
dexter-	1/1	Running	0	8m22s	10.	180		us-east-2.compute.internal	<none></none>	<none></none>

Now I deleted the deployment dexter and found that two new scaled nodes had been deleted as shown in the screenshot attached below.

```
[root@ ~]# kubectl delete -f deployment.yaml
deployment.apps "dexter" deleted
[root@ ~]# kubectl get pods --watch
^C[root@ ~]#
```

After unneed for 10 minutes the newly added nodes were deleted as can be seen from the logs of cluster-autoscaler pod present in the namespace kube-system.

```
| Foot@ | Testic_autoscaler_go:589 | Scale down status: lastScalephrise:205-08-16 11:01:36.69431923 +0000 UTC m=-487.61918864 lastScaleDownPailTime=2025-08-16 09:55:33.20551210 +0000 UTC m=-3593.877000553 scaleDownPailTime=2025-08-16 09:55:33.20551210 +0000 UTC m=-3593.877000553 scaleDownPailTime=2025-08-15 (scaleDownPailTime=2025-08-15 (scaleDownPailTime=2025-08-
```

## 2. Second Section

#### **Karpenter**

Karpenter is an open source kubernetes cluster autoscaler. Karpenter is 10 times faster than cluster autoscaler. I had created the EKS Cluster using the terraform script present in my GitHub Repo <a href="https://github.com/singhritesh85/Autoscale.git">https://github.com/singhritesh85/Autoscale.git</a> at the path Autoscale/Autoscale-EKS/Autoscale-EKS-Cluster-using-Karpenter and edited the configmap aws-auth present at the namespace kube-system then created the Karpenter controller using the helm run the manifests file karpenter.yaml as mentioned below.

```
~]# kubectl get nodes --watch
[root@
                                                             AGE
NAME
                                           STATUS
                                                    ROLES
                                                                   VERSION
ip-10-
          -12.us-east-2.compute.internal
                                           Ready
                                                             12m
                                                                   v1.30.4-eks-a
                                                    <none>
ip-10-
          -158.us-east-2.compute.internal
                                          Ready
                                                    <none>
                                                             12m
                                                                 v1.30.4-eks-a
```

Then create the deployment with 40 pods and found that pods are in pending state due to unavailability of compute resources. Which you can check with the command **kubectl describe pod <pod-name> -n namespace**.

```
~]# vim deployment.yaml
[root@
                    ~]# kubectl apply -f deployment.yaml
[root@
deployment.apps/dexter created
[root@
                    ~]# cat deployment.yaml
apiVersion: apps/v1
kind: Deployment
metadata:
 name: dexter
spec:
 selector:
   matchLabels:
     app: nginx
 replicas: 40
  template:
    metadata:
     labels:
       app: nginx
    spec:
      containers:
      - name: nginx
       image: nginx
       ports:
        - containerPort: 80
              ~]# kubectl get pods
[root@:
NAME
                          READY
                                  STATUS
                                            RESTARTS
                                                       AGE
                          0/1
                                  Pending
dexter-
                                                       19s
dexter-
                          0/1
                                  Pending
                                                       18s
                                  Pending
dexter-
                          0/1
                                                       18s
dexter-
                          1/1
                                  Running
                                           0
                                                       19s
dexter-
                                  Pending 0
                          0/1
                                                       18s
                                          0
dexter-
                                  Running
                                                       19s
                          1/1
                                                        19s
dexter-
                          1/1
                                  Running
                                            0
dexter-
                          0/1
                                  Pending
                                            0
                                                       18s
dexter-
                          1/1
                                  Running
                                            0
                                                       19s
dexter-
                          0/1
                                  Pending
                                           9
                                                       18s
                          1/1
                                  Running
                                                       19s
dexter-
dexter-
                          1/1
                                  Running 0
                                                       19s
```

[root@ ~]# kubectl edit cm aws-auth -n kube-system

```
# Please edit the object below. Lines beginning with a '#' will be ignored,
# and an empty file will abort the edit. If an error occurs while saving this file will be
# reopened with the relevant failures.
apiVersion: v1
data:
 mapRoles:
   - groups:
     - system:bootstrappers
     - system:nodes
                                 6:role/eks-nodegroup-role-dev
    rolearn: arn:aws:iam::0
     username: system:node:{{EC2PrivateDNSName}}
   - groups:
      system:bootstrappers
     - system:nodes
     rolearn: arn:aws:iam::0
                                   6:role/karpenter-eks-noderole
     username: system:node:{{EC2PrivateDNSName}}
kind: ConfigMap
metadata:
 creationTimestamp: "2025-
 name: aws-auth
 namespace: kube-system
 resourceVersion: "___"
 uid: b
-- INSERT (paste) --
```

Install the karpenter controller using the helm as shown in the screenshot attached below.

helm upgrade --install karpenter oci://public.ecr.aws/karpenter/karpenter --version "1.3.2" -namespace "karpenter" --create-namespace --set "settings.clusterName=eks-demo-cluster-dev" -set "serviceAccount.annotations.eks\.amazonaws\.com/rolearn=arn:aws:iam::02XXXXXXXXXXXC:role/karpenter-controller-role" --set
controller.resources.requests.cpu=1 --set controller.resources.requests.memory=1Gi --set
controller.resources.limits.cpu=1 --set controller.resources.limits.memory=1Gi

The two pods present in the namespace karpenter are for karpenter controller and then run the manifests file karpenter.yaml as shown in the screenshot attached below.

```
[root@______terraform-eks-withaddons-karpenter]# kubectl apply -f karpenter.yaml nodepool.karpenter.sh/general-purpose created ec2nodeclass.karpenter.k8s.aws/default created
```

For your reference I kept this yaml manifests file in GitHub Repo <a href="https://github.com/singhritesh85/Autoscale.git">https://github.com/singhritesh85/Autoscale.git</a> at the path Autoscale/Autoscale-EKS/Autoscale-EKS-Cluster-using-Karpenter/terraform-eks-withaddons-karpenter. You can edit this yaml manifests file depending on your project then apply.

```
terraform-eks-withaddons-karpenter]# cat karpenter.yaml
# This example NodePool will provision general purpose instances
apiVersion: karpenter.sh/v1
kind: NodePool
metadata:
 name: general-purpose
   kubernetes.io/description: "General purpose NodePool for generic workloads"
 template:
   spec:
     requirements:
         key: kubernetes.io/arch
         operator: In
values: ["amd64"]
         key: kubernetes.io/os
         operator: In
values: ["linux"]
        - key: karpenter.sh/capacity-type
         operator: In
                                  #["spot"] ### You can select on-demand or spot instance depending on your requirement.
        - key: karpenter.k8s.aws/instance-category
         operator: In
         values: ["t"]
                                                                            #["c", "m", "r"]
                          # Interested to launch t series instance
        - key: karpenter.k8s.aws/instance-generation
         operator: Gt
         values: ["2"]
                         # Instances launched will be greated than 2
     nodeClassRef:
       group: karpenter.k8s.aws
       kind: EC2NodeClass
       name: default
apiVersion: karpenter.k8s.aws/v1
kind: EC2NodeClass
metadata:
 name: default
 annotations
   kubernetes.io/description: "General purpose EC2NodeClass for running Amazon Linux 2 nodes"
 role: "karpenter-eks-noderole" # replace with your karpenter noderole
 subnetSelectorTerms:
   - tags:
       karpenter.sh/discovery: "eks-demo-cluster-dev" # replace with your cluster name
 securityGroupSelectorTerms:
   - tags:
       karpenter.sh/discovery: "eks-demo-cluster-dev" # replace with your cluster name
 amiSelectorTerms:
     alias: al2023@latest # Amazon Linux 2023
                    terraform-eks-withaddons-karpenter]#
```

Then I checked the logs of the karpenter controller pod present in the namespace karpenter as shown in the screenshot attached below.

```
ts":{"cpu":"388m","memory":"376Hi","pods":"26"},"instance-types":"t3.2xlarge, t3.1arge, t3.2xlarge, t3a.2xlarge, t3a.1arge and 1 other(s)"}

("level":"INFO", "time":"2825-83-1718:15:32.8952", "logge":"controller", "message": "launched nodeclaim", "commit":"1 6", "controller":"nodeclaim.lifecycle", "controllerGroup:"Karpenter.sh", "controllerKind": Nodeclaim":("name":"general-purpose-we "v"), "namespace":"", "name":"general-purpose-we 76", "necor citaTO:"0 5", "provider-id":"ans::///ivs-asst-2/i-0 7", "instance-type":"t3.1arge", "lone"; "us-asst-2a", capacity -type":"on-demand", "allocatable":("cpu":"1938m", "ephemeral-storage":"170i", "memory":"6837Hi", "pods":"35")}

[root@ip-172-31-3-49 main]# [root@ip-172-31-3-49 main]# [whectl logs -f Karpenter-8]
```

Finally, I found a new node had been created which joined the cluster as shown in the screenshot attached below.

```
[root@
                        ]# kubectl get nodes
NAME
                                          STATUS
                                                   ROLES
                                                            AGE
                                                                 VERSTON
                                                   <none>
                                                           93m v1.30.4-eks-a737599
ip-10-
                                          Ready
          -65.us-east-2.compute.internal
ip-10-
         -148.us-east-2.compute.internal
                                          Ready
                                                           33m
                                                                 v1.30.4-eks-a737599
                                                   (none)
         -58.us-east-2.compute.internal
                                                   <none> 94s v1.30.9-eks-5d632ec
ip-10-
                                          Ready
```

You can provide your choice as per your project requirement in the file karpenter.yaml like which series instance type you want to get launched, whether you want on-demand or spot instances and other stuff.

I would like to mention here in industry at most of the cases the kubernetes cluster is available with the cluster-autoscaler and if you want to migrate to karpenter from cluster-autoscaler then do the below changes in your existing EKS Cluster.

- 1. Delete the yaml manifests for cluster-autoscaler using the command kubectl delete -f cluster-autoscaler.yaml.
- 2. Add the tags to subnets and security groups of the Nodes which were the part that EKS Cluster.
- 3. Add the IAM Role and Policy for Karpenter Controller and Karpenter Nodes Roles and Karpenter Nodes Policy (This changes you can perform in your existing terraform script).
- 4. Edit the configmap aws-auth present in the namespace kube-system.
- 5. Install the karpenter controller using the helm.
- 6. Do the changes in karpenter.yaml as per your project requirement like which series of instance you want to be get launched, on-demand/spot instances and other stuff then apply the manifests file using the command kubectl apply -f karpenter.yaml.

```
[root@ ~]# cd terraform-eks-withaddons-cluster-autoscaler/
[root@ terraform-eks-withaddons-cluster-autoscaler]# kubectl delete -f cluster-autoscaler.yaml
serviceaccount "cluster-autoscaler" deleted
clusterrole.rbac.authorization.k8s.io "cluster-autoscaler" deleted
role.rbac.authorization.k8s.io "cluster-autoscaler" deleted
clusterrolebinding.rbac.authorization.k8s.io "cluster-autoscaler" deleted
rolebinding.rbac.authorization.k8s.io "cluster-autoscaler" deleted
deployment.apps "cluster-autoscaler" deleted
```

Below changes had been done in vpc.tf file.

```
resource "aws_subnet" "public_subnet" {
 count = "${length(data.aws_availability_zones.azs.names)}"
 vpc_id
         = "${aws_vpc.test_vpc.id}"
 availability_zone = "${element(data.aws_availability_zones.azs.names,count.index)}"
 cidr_block = "${element(var.public_subnet_cidr,count.index)}"
 map_public_ip_on_launch = true
 tags = {
   Name = "PublicSubnet-${var.env}-${count.index+1}"
   Environment = var.env
                              ##"${terraform.workspace}"
   "karpenter.sh/discovery" = "${var.eks_cluster}-${var.env}"
resource "aws_subnet" "private_subnet" {
 count = "${length(data.aws_availability_zones.azs.names)}"
                                                                 ##"${length
         = "${aws_vpc.test_vpc.id}"
 availability_zone = "${element(data.aws_availability_zones.azs.names,count.index)}"
 cidr_block = "${element(var.private_subnet_cidr,count.index)}"
 tags = {
   Name = "PrivateSubnet-${var.env}-${count.index+1}"
                                  ##"${terraform.workspace}"
   Environment = var.env
   "karpenter.sh/discovery" = "${var.eks_cluster}-${var.env}"
```

Open the file cluster.tf and go to the end and edit with content as written below.

```
# Tag the EKS Cluster Security Group
resource "aws_ec2_tag" "cluster_security_group" {
resource_id = aws_eks_cluster.eksdemo.vpc_config[0].cluster_security_group_id
     = "karpenter.sh/discovery"
key
     = "${var.eks_cluster}-${var.env}"
value
depends_on = [aws_eks_node_group.eksnode]
}
# Karpenter Controller
data "aws_iam_policy_document" "karpenter_assume_role_policy" {
statement {
 actions = ["sts:AssumeRoleWithWebIdentity"]
 effect = "Allow"
 condition {
 test = "StringEquals"
 variable = "${replace(aws_iam_openid_connect_provider.eksopidc.url, "https://", "")}:sub"
 values = ["system:serviceaccount:karpenter:karpenter"]
```

}

principals {

```
identifiers = [aws_iam_openid_connect_provider.eksopidc.arn]
   type
           = "Federated"
  }
}
resource "aws_iam_role" "karpenter_controller_role" {
 assume_role_policy = data.aws_iam_policy_document.karpenter_assume_role_policy.json
              = "karpenter-controller-role"
 name
}
resource "aws_iam_role_policy_attachment" "karpenter_controller_policy_attach" {
 role
        = aws_iam_role.karpenter_controller_role.name
 policy_arn = aws_iam_policy.karpenter_controller_policy.arn
}
resource "aws_iam_policy" "karpenter_controller_policy" {
 name = "karpenter-controller-policy"
 policy = jsonencode(
   "Version": "2012-10-17",
   "Statement":[
     "Sid": "VisualEditor0",
     "Effect": "Allow",
     "Action": "ec2:TerminateInstances",
     "Resource":
"arn:aws:ec2:${data.aws_region.reg.id}:${data.aws_caller_identity.G_Duty.account_id}:instance/*",
     "Condition": {
      "ForAnyValue:StringLike" : {
```

```
"ec2:DeleteLaunchTemplate",
     "ec2:CreateLaunchTemplate",
     "ec2:CreateTags"
    ],
    "Resource":
template/*"
   },
    "Sid": "VisualEditor3",
    "Effect": "Allow",
    "Action" : [
     "iam:PassRole",
     "ssm:GetParameter"
    ],
    "Resource" : [
     "arn:aws:iam::${data.aws_caller_identity.G_Duty.account_id}:role/karpenter-node*",
     "arn:aws:ssm:${data.aws_region.reg.id}:*:parameter/*"
    ]
   },
    "Sid": "karpenterInstanceProfile",
    "Effect": "Allow",
    "Action" : [
     "iam:CreateInstanceProfile",
     "iam:TagInstanceProfile",
     "iam:AddRoleToInstanceProfile",
     "iam:RemoveRoleFromInstanceProfile",
     "iam:DeleteInstanceProfile"
    ],
    "Resource":[
```

```
"arn:aws:iam::${data.aws_caller_identity.G_Duty.account_id}:instance-profile/eks-demo-
cluster-dev*"
   ]
}
####################################
# Nodes Spinned-up using Karpenter
# IAM Role which allows Karpenter to spin up nodes
data "aws_iam_policy_document" "karpenter_node_iam_role_policy" {
statement {
 actions = ["sts:AssumeRole"]
 effect = "Allow"
 principals {
  identifiers = ["ec2.amazonaws.com"]
  type
        = "Service"
resource "aws_iam_role" "karpenter_node_iam_role" {
assume_role_policy = data.aws_iam_policy_document.karpenter_node_iam_role_policy.json
          = "karpenter-eks-noderole"
name
```

```
}
resource "aws_iam_role_policy_attachment" "karpenter_workernode" {
 role
        = aws_iam_role.karpenter_node_iam_role.name
 policy_arn = "arn:aws:iam::aws:policy/AmazonEKSWorkerNodePolicy"
}
resource "aws_iam_role_policy_attachment" "karpenter_cni" {
 role
        = aws_iam_role.karpenter_node_iam_role.name
 policy_arn = "arn:aws:iam::aws:policy/AmazonEKS_CNI_Policy"
}
resource "aws_iam_role_policy_attachment" "karpenter_ecr"
 role
        = aws_iam_role.karpenter_node_iam_role.name
 policy_arn = "arn:aws:iam::aws:policy/AmazonEC2ContainerRegistryReadOnly"
}
resource "aws_iam_role_policy_attachment" "karpenter_ssm" {
        = aws_iam_role.karpenter_node_iam_role.name
 policy_arn = "arn:aws:iam::aws:policy/AmazonSSMManagedInstanceCore"
}
resource "aws_iam_instance_profile" "karpenter_instance_profile" {
 name = "karpenter-eks-instance-profile"
 role = aws_iam_role.karpenter_node_iam_role.name
}
```

```
"ec2:ResourceTag/Name" : "*karpenter*"
  }
}
},
 "Sid": "VisualEditor1",
 "Effect": "Allow",
 "Action" : [
  "eks:DescribeCluster",
  "ec2:DescribeSpotPriceHistory",
  "ec2:DescribeImages",
  "ec2:DescribeInstances",
  "ec2:DescribeInstanceTypeOfferings",
  "ec2:DescribeAvailabilityZones",
  "ec2:DescribeLaunchTemplates",
  "ec2:DescribeInstanceTypes",
  "ec2:RunInstances",
  "ec2:CreateFleet",
  "ec2:DescribeSubnets",
  "ec2:DescribeSecurityGroups",
  "ec2:CreateTags",
  "pricing:GetProducts",
  "iam:GetInstanceProfile",
  "iam:PassRole"
 "Resource": "*"
},
 "Sid": "VisualEditor2",
 "Effect": "Allow",
 "Action" : [
```

```
"ec2:DeleteLaunchTemplate",
     "ec2:CreateLaunchTemplate",
     "ec2:CreateTags"
    ],
    "Resource":
template/*"
   },
    "Sid": "VisualEditor3",
    "Effect": "Allow",
    "Action" : [
     "iam:PassRole",
     "ssm:GetParameter"
    ],
    "Resource" : [
     "arn:aws:iam::${data.aws_caller_identity.G_Duty.account_id}:role/karpenter-node*",
     "arn:aws:ssm:${data.aws_region.reg.id}:*:parameter/*"
    ]
   },
    "Sid": "karpenterInstanceProfile",
    "Effect": "Allow",
    "Action" : [
     "iam:CreateInstanceProfile",
     "iam:TagInstanceProfile",
     "iam:AddRoleToInstanceProfile",
     "iam:RemoveRoleFromInstanceProfile",
     "iam:DeleteInstanceProfile"
    ],
    "Resource":[
```

```
"arn:aws:iam::${data.aws_caller_identity.G_Duty.account_id}:instance-profile/eks-demo-
cluster-dev*"
   ]
   }
}
#############################
# Nodes Spinned-up using Karpenter
# IAM Role which allows Karpenter to spin up nodes
data "aws_iam_policy_document" "karpenter_node_iam_role_policy" {
statement {
 actions = ["sts:AssumeRole"]
 effect = "Allow"
 principals {
  identifiers = ["ec2.amazonaws.com"]
        = "Service"
  type
resource "aws_iam_role" "karpenter_node_iam_role" {
assume_role_policy = data.aws_iam_policy_document.karpenter_node_iam_role_policy.json
          = "karpenter-eks-noderole"
name
```

```
}
resource "aws_iam_role_policy_attachment" "karpenter_workernode" {
 role
        = aws_iam_role.karpenter_node_iam_role.name
 policy_arn = "arn:aws:iam::aws:policy/AmazonEKSWorkerNodePolicy"
}
resource "aws_iam_role_policy_attachment" "karpenter_cni" {
 role
        = aws_iam_role.karpenter_node_iam_role.name
 policy_arn = "arn:aws:iam::aws:policy/AmazonEKS_CNI_Policy"
}
resource "aws_iam_role_policy_attachment" "karpenter_ecr"
 role
        = aws_iam_role.karpenter_node_iam_role.name
 policy_arn = "arn:aws:iam::aws:policy/AmazonEC2ContainerRegistryReadOnly"
}
resource "aws_iam_role_policy_attachment" "karpenter_ssm" {
 role
        = aws_iam_role.karpenter_node_iam_role.name
 policy_arn = "arn:aws:iam::aws:policy/AmazonSSMManagedInstanceCore"
}
resource "aws_iam_instance_profile" "karpenter_instance_profile" {
 name = "karpenter-eks-instance-profile"
 role = aws_iam_role.karpenter_node_iam_role.name
}
```

I changed the min\_size, max\_size and desired\_capacity and kept them same. Take the desired count such as it could bear all the pods required for the exiting projects running in this cluster. For the current scenario I kept it as 2.

```
scaling_config {
  desired_size = 2
  max_size = 2
  min_size = 2
}
```

dexter-6cfb64b7c5-574pk 1/1

dexter-6cfb64b7c5-57sbv 0/1

dexter-6cfb64b7c5-77d8q 1/1

dexter-6cfb64b7c5-7fqj5 1/1 Running

```
main]# terraform apply -auto-approve
[root@
[root@
                   ~]# cat deployment.yaml
apiVersion: apps/v1
kind: Deployment
metadata:
 name: dexter
spec:
  selector:
    matchLabels:
      app: nginx
  replicas: 30
  template:
    metadata:
      labels:
        app: nginx
    spec:
      containers:
      - name: nginx
        image: nginx
        ports:
        containerPort: 80
              ~]# kubectl apply -f deployment.yaml
[root@
[root@ip-172-31-3-49 ~]# kubectl get pods --watch
                       READY
                               STATUS
                                                 RESTARTS
                                                           AGE
                       1/1
dexter-6cfb64b7c5-2tdz2
                               Running
                                                            15s
dexter-6cfb64b7c5-4bwah
                       1/1
                               Running
                                                 0
                                                            15s
dexter-6cfb64b7c5-4sbk7 0/1
                               Pending
                                                 0
                                                            15s
```

Running

Pending

Running

15s

15s

15s

15s

0

0

0

Pod Disruption Budget must had applied with the existing pods for your projects. Otherwise, if the application becomes unavailable then DevOps team which Email-ID was provided to the notification for Synthetic Monitoring (Unavailability for endpoint URL) will get notification. So please do this activity during non-production hours or take-off the email-id from the Alerting (However it is not suggestable to remove the DevOps team Email-ID from the Alerting) and after the activity add back to Alerting.



```
# This example NodePool will provision general purpose instances
apiVersion: karpenter.sh/v1
kind: NodePool
metadata:
name: general-purpose
annotations:
  kubernetes.io/description: "General purpose NodePool for generic workloads"
spec:
template:
  spec:
   requirements:
    - key: kubernetes.io/arch
     operator: In
     values: ["amd64"]
    - key: kubernetes.io/os
     operator: In
     values: ["linux"]
    - key: karpenter.sh/capacity-type
     operator: In
     values: ["on-demand"] #["spot"] ### You can select on-demand or spot instance depending
on your requirement.
    - key: karpenter.k8s.aws/instance-category
     operator: In
     values: ["t"] # Interested to launch t series instance #["c", "m", "r"]
    - key: karpenter.k8s.aws/instance-generation
     operator: Gt
     values: ["2"] # Instances launched will be greated than 2
   nodeClassRef:
    group: karpenter.k8s.aws
    kind: EC2NodeClass
```

name: default

---

apiVersion: karpenter.k8s.aws/v1

kind: EC2NodeClass

metadata:

name: default

annotations:

kubernetes.io/description: "General purpose EC2NodeClass for running Amazon Linux 2 nodes" spec:

role: "karpenter-eks-noderole" # replace with your karpenter noderole subnetSelectorTerms:

- tags:

karpenter.sh/discovery: "eks-demo-cluster-dev" # replace with your cluster name securityGroupSelectorTerms:

- tags:

karpenter.sh/discovery: "eks-demo-cluster-dev" # replace with your cluster name amiSelectorTerms:

- alias: al2023@latest # Amazon Linux 2023

Checked the logs of karpenter controller pod and found a node had been created as shown in the screenshot attached below.

[root@	#	kubect]	l get pods		
NAME		READY	STATUS	RESTARTS	AGE
dexter-6	2	1/1	Running	0	m
dexter-6	h	1/1	Running	0	m
dexter-6	7	1/1	Running	0	m
dexter-6	k	1/1	Running	0	m
dexter-6	V	1/1	Running	0	m
dexter-6	q	1/1	Running	0	m
dexter-6	5	1/1	Running	0	m
dexter-6	9 5 c 5	1/1	Running	0	m
dexter-6	5	1/1	Running	0	m
dexter-6	W	1/1	Running	0	m
dexter-6	h	1/1	Running	0	m
dexter-6	h	1/1	Running	0	m
dexter-6	8	1/1	Running	0	m
dexter-6	q t	1/1	Running	0	m
dexter-6	t	1/1	Running	0	m
dexter-6	q	1/1	Running	0	m
dexter-6	W	1/1	Running	0	m
dexter-6	5 5 X	1/1	Running	0	m
dexter-6	5	1/1	Running	0	m
dexter-6	X	1/1	Running	0	m
dexter-6	s	1/1	Running	0	m
dexter-6	s z n	1/1	Running	0	m
dexter-6	n	1/1	Running	0	m
dexter-6	h	1/1	Running	0	m
dexter-6	V	1/1	Running	0	m
dexter-6	f	1/1	Running	0	m
dexter-6	f 6 4	1/1	Running	0	m
dexter-6	4	1/1	Running	0	m
dexter-6	С	1/1	Running	0	m
dexter-6	n	1/1	Running	0	m
[root@ip-172-31-3-49 ~]# kubec	tl go	at nodes			
NAME	cr ge		TATUS ROLES	AGE VERSION	N
ip-10140.us-east-2.compu	te.in	nternal R	eady <none></none>	v1.30.4	4-eks-a 9

```
-139.us-east-2.compute.internal
                                          Ready
                                                  <none>
                                                                  v1.30.4-eks-a
ip-10-
        -79.us-east-2.compute.internal Ready
                                                                  v1.30.9-eks-5
                                                  <none>
```

The above attached screenshot verifies the autoscaling had been done and migrated from clusterautoscaler to karpenter

### 3. Third Section

### **Cluster-Autoscaler in AKS Cluster**

I used the terraform script to created the AKS Cluster and enabled the cluster autoscaler in default node pool (system node pool) and user node pool by keeping **auto\_scaling\_enabled = true**. As shown in the screen shot attached below, I kept the min\_count and max\_count as 1 and 4 respectively, the node\_count which indicates the desired node count kept as 1. The node\_count cannot be less than min\_count and cannot be less than max\_count.

```
min_count ≤ node_count ≤ max_count

auto_scaling_enabled = true
max_count = 4
node_count = 1
min_count = 1
```

Initially just after creation of the AKS Cluster there were 2 nodes one for system pool and one for agent pool.

```
[root@
                      ~]# kubectl get nodes
NAME
                                     STATUS
                                               ROLES
                                                        AGE
                                                               VERSION
aks-agentpool-2
                     4-vmss000000
                                                              v1.29.0
                                     Readv
                                               <none>
aks-userpool-2
                     4-vmss000000
                                     Ready
                                               <none>
                                                               v1.29.0
```

You can check the logs of cluster autoscaling using the command **kubectl get events -n kube-system** or **kubectl describe cm cluster-autoscaler-status -n kube-system**.

For this demonstration I created a deployment with nginx image and kept the number of desired pods as 250.

```
[root@]
                       ~]# cat deployment.yaml
apiVersion: apps/v1
kind: Deployment
metadata:
 name: dexter
spec:
  selector:
   matchLabels:
      app: nginx
 replicas: 250
 template:
   metadata:
      labels:
        app: nginx
    spec:
      containers:
      - name: nginx
        image: nginx
        ports:
        - containerPort: 80
                  ~]# kubectl apply -f deployment.yaml
```

[root@	~]#	# kubect:	l get pods	watch	
NAME		READY	STATUS	RESTARTS	AGE
dexter-7	)	1/1	Running	0	56s
dexter-7	N	0/1	Pending	0	19s
dexter-7	2	0/1	Pending	0	16s
dexter-7	)	1/1	Running	0	55s
dexter-7	j	0/1	Pending	0	20s
dexter-7	g	0/1	Pending	0	16s
dexter-7	F	1/1	Running	0	55s
dexter-7	g	1/1	Running	0	55s
dexter-7	)	1/1	Running	0	21s
dexter-7	F	0/1	Pending	0	16s
dexter-7	7	0/1	Pending	0	16s
dexter-7	L	1/1	Running	0	22s
dexter-7 r	า	1/1	Running	0	91s
dexter-7	K	1/1	Running	0	21s
dexter-7	5	0/1	Pending	0	19s
dexter-7	l	1/1	Running	0	92s
dexter-7	9	0/1	Pending	0	16s
dexter-7	9	0/1	Pending	0	16s

As there were lack of compute resources due to which pods could not be scheduled on the nodes and are in pending state which can be seen in the screenshot attached above, you can also check the same using the command **kubectl describe pod <pod-name> -n namespace**.

Now the cluster-autoscaler in AKS will spin new nodes on which the pods can be scheduled and I found that cluster-autoscaler spun 2 new nodes so that it could schedule the pods to the nodes.

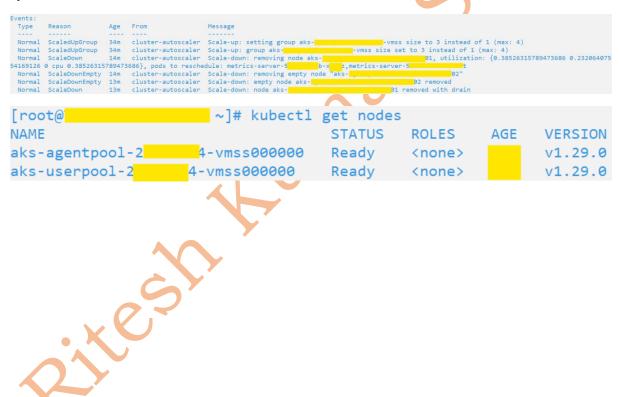
[root@	_~]# kub∈	ectl	get no	des			
NAME			STATU	S	ROLES	AGE	VERSION
aks-		00	Ready		<none></none>		v1.29.0
aks-		01	Ready		<none></none>		v1.29.0
aks-		02	Ready		<none></none>		v1.29.0
aks-		0	Ready		<none></none>		v1.29.0
aks-		O	Ready		(Hone)		V1.29.0
[root@	~]# kubec	tl ge	t pods				
NAME	READY	STA		RES1	TARTS	AGE	
dexter-7	1/1	Run	ning	0			
dexter-7	1/1		ning	0			
dexter-7	1/1		ning	0			
dexter-7	1/1		ning	0			
dexter-7	1/1		ning	0			
dexter-7	1/1		ning	0			
dexter-7	1/1		ning	0			
dexter-7	1/1		ning	0			
dexter-7	1/1		ning	0			
dexter-7	1/1		ning	0			
dexter-7	1/1		ning	0			
dexter-7	1/1		ning	0			
dexter-7	1/1		ning	0			
dexter-7			_				
dexcer-/	1/1	Kun	ning	0			

I checked the logs using the command **kubectl describe cm cluster-autoscaler-status -n kube-system** as shown below.

Finally, I deleted the deployment dexter using the command as shown in the screenshot attached below and found that no pod is running now.

```
[root@ ~]# kubectl delete -f deployment.yaml |
[root@ ~]# kubectl get pods
No resources found in default namespace.
```

Then checked the logs using the command **kubectl describe cm cluster-autoscaler-status -n kube-system** as shown below and found that nodes are scaled down to the initial state.



**GitHub Repository:** - https://github.com/singhritesh85/Autoscale.git

#### References: -

- 1. <a href="https://github.com/kubernetes/autoscaler/blob/master/cluster-autoscaler-autoscaler-autodiscover.yaml">https://github.com/kubernetes/autoscaler/blob/master/cluster-autoscaler-autodiscover.yaml</a>
- 2. https://karpenter.sh/docs/getting-started/getting-started-with-karpenter/
- 3. <a href="https://github.com/aws/karpenter-provider-aws/blob/main/examples/v1/general-purpose.yaml">https://github.com/aws/karpenter-provider-aws/blob/main/examples/v1/general-purpose.yaml</a>