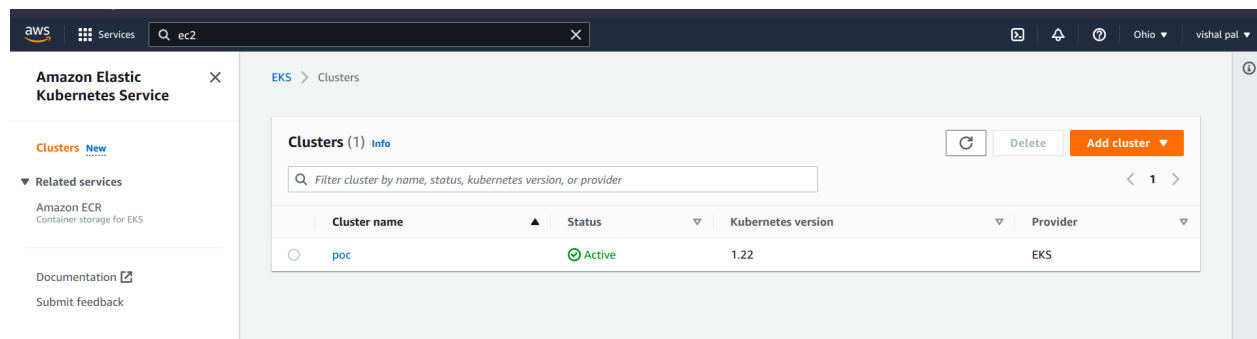# USING SPOT INSTANCE FOR DEPLOYMENT WITH NO DOWNTIME DURING RECLAMATION OF SPOT NODE USING NODE TERMINATION HANDLER.

## CREATING EKS CLUSTER

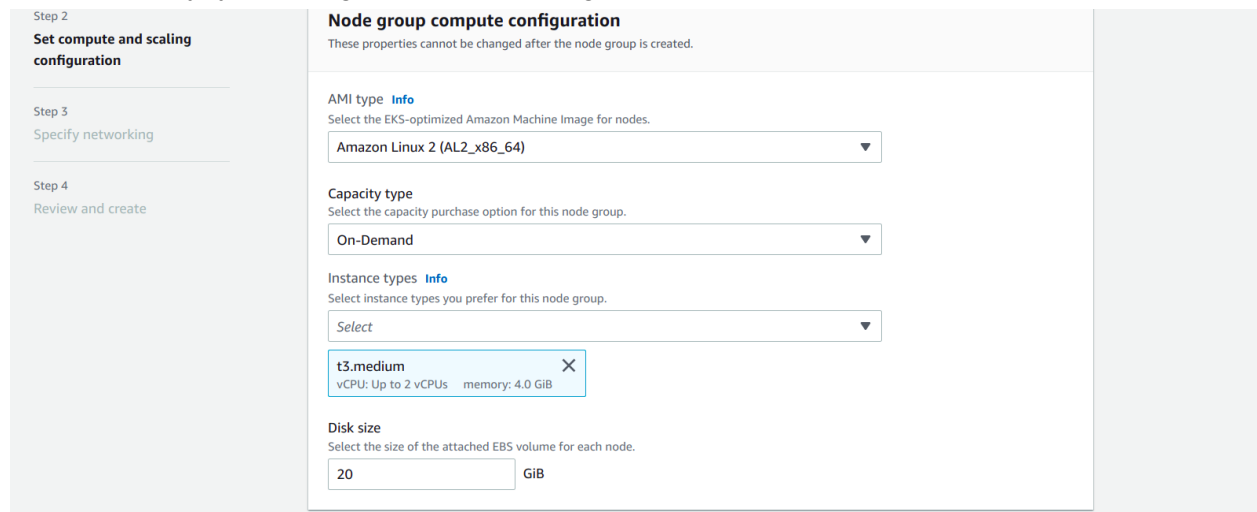We can refer this document to create a kubernetes cluster.

https://docs.aws.amazon.com/eks/latest/userguide/create-cluster.html.



## CREATING A ON_DEMAND NODE GROUP

We can refer to this document to create a aws node group.

https://docs.aws.amazon.com/eks/latest/userguide/create-managed-node-group.html

In the Capacity type during creation of node group we will select ON_DEMAND



After creation of an on_demand node group we will create a different node group for instance.

We will choose the minimum available node as 1

            Maximum available node as 1

            Desired nodes as 1

As we are only going to scale SPOT nodes.


CREATING A ON_SPOT NODE GROUP

We can refer to this document to create a spot aws node group.

https://docs.aws.amazon.com/eks/latest/userguide/create-managed-node-group.html

During creation of node group we will select capacity type as spot instance



We make sure  If we selected Spot for Capacity type, then we recommend specifying multiple instance types to enhance availability.

We will choose the minimum available node as 1

            Maximum available node as 4

            Desired nodes as 1


Make sure to add a label to this node group.

type = SPOT


After creation of node groups we will have one available node for each node group.

## CREATING A DEPLOYMENT AND A SERVICE

We will create a deployment deployment.yaml file and apply it to launch the pods
$ Kubectl apply -f deployment.yaml

Next we will create a service file for the above deployment
$ Kubectl apply -f service.yaml

We will serve our deployment using ingress


## INSTALLING INGRESS USING HELM

Enter the following command to add the Helm ingress-nginx repo:
$ helm repo add ingress-nginx https://kubernetes.github.io/ingress-nginx

Command to create a namespace for ingress
```
$ kubectl create namespace nginx-ingress-sample
```

```
Command to deploy ingress resources using helm
$ helm install my-nginx ingress-nginx/ingress-nginx \
--namespace nginx-ingress-sample \
--set controller.metrics.enabled=true \
--set-string
controller.metrics.service.annotations."prometheus\.io/port"="10254" \
--set-string
controller.metrics.service.annotations."prometheus\.io/scrape"="true"
```

Now we will create a ingress file for our service

```
unthinkable-lap-0200@PG02PVPL:~/Desktop/aws$ cat myingress
apiVersion: networking.k8s.io/v1
kind: Ingress
metadata:
  name: minimal-ingress
  annotations:
    nginx.ingress.kubernetes.io/rewrite-target: /
spec:
  ingressClassName: nginx
  #host:
  rules:
  - http:
      paths:
      - path: /pop
        pathType: Prefix
        backend:
          service:
            name: svc-hello
            port:
              number: 8080
```

Now we will apply the file
$ kubectl apply -f <ingress-file-name>

And now we will get our ingress load balancer url using command:
$ kubectl get ingress

```
unthinkable-lap-0200@PG02PVPL:~/Desktop/aws$ kubectl get ingress
NAME              CLASS   HOSTS   ADDRESS                                                                          PORTS   AGE
minimal-ingress   nginx   *       a29800e773c144b6cb04cadf8108458f-160032197.us-east-2.elb.amazonaws.com   80      2d9h
unthinkable-lap-0200@PG02PVPL:~/Desktop/aws$
```

## INSTALLING KUBERNETES METRICS SERVER
https://docs.aws.amazon.com/eks/latest/userguide/metrics-server.html

Metrics Server is a scalable, efficient source of container resource metrics for
Kubernetes built-in autoscaling pipelines.
**Use Case**
- CPU/Memory based horizontal autoscaling (learn more about Horizontal Autoscaling)
- Automatically adjusting/suggesting resources needed by containers (learn more about Vertical Autoscaling)

We can use this commands:

$ kubectl apply -f
https://github.com/kubernetes-sigs/metrics-server/releases/latest/download/components.yaml

To get metrics server deployment

$ kubectl get deployment metrics-server -n kube-system


INSTALLING HORIZONTAL POD AUTOSCALER

https://docs.aws.amazon.com/eks/latest/userguide/horizontal-pod-autoscaler.html

We can directly use this commands :
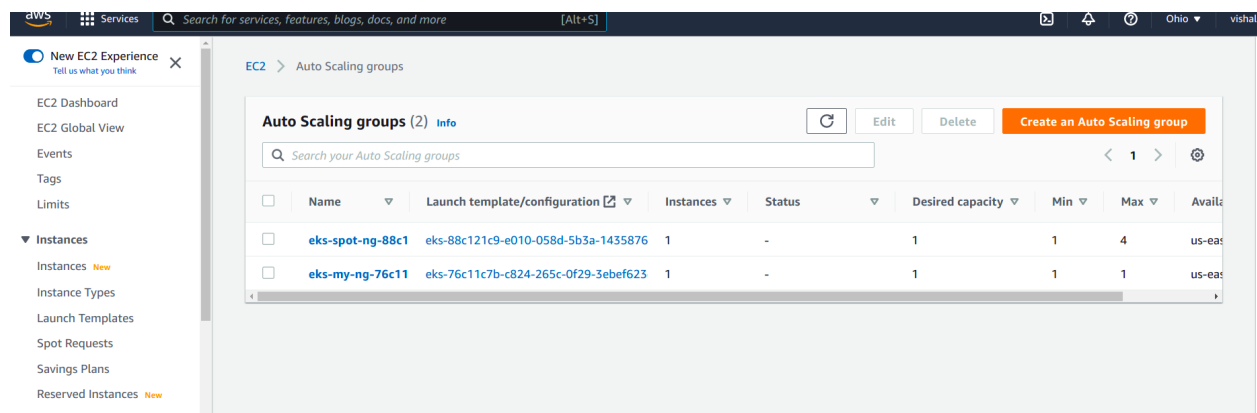$ kubectl autoscale deployment <deployment-name> --cpu-percent=50 --min=1 --max=10



INSTALLING CLUSTER AUTOSCALER ON CLUSTER

The Kubernetes Cluster Autoscaler automatically adjusts the number of nodes in your cluster
when pods fail or are rescheduled onto other nodes. The Cluster Autoscaler is typically installed
as a Deployment in your cluster.

We can refer to this document to install an autoscaler in our eks cluster:
https://docs.aws.amazon.com/eks/latest/userguide/autoscaling.html


After installing cluster autoscaler it wil create a autoscaling group for both out our node group in
AWS autoscaling groups console



We will dynamically scale the spot node group
For creating a dynamic scaling policy
Go to auto scaling groups →Automatic scaling → Dynamic scaling policies →create a Dynamic
Autoscaling policy.

We will choose cpu usage to our desired usage after that we want to scale out nodes.
And warmup time for nodes .

Make Sure scale in is enabled.

INSTALLING NODE TERMINATION HANDLER ON CLUSTER

To install node termination handler we will use these command

$ helm repo add eks https://aws.github.io/eks-charts

$ helm install aws-node-termination-handler \
        --namespace kube-system \
        --version 0.15.4 \
        --set nodeSelector.type= SPOT\
        eks/aws-node-termination-handler

To manage the Availability of our website we will create two deployments with a single service to handle no downtime during the reclamation of nodes by aws.

To make sure there is at least one pod available on each node we will use selectors to deploy on each node .

We can generate load on out load balancer endpoint using

$ kubectl run -i \
   --tty load-generator \
   --rm --image=busybox \
   --restart=Never \
   -- /bin/sh -c "while sleep 0.01; do wget -q -O- <out-end-point>; done"

To get HPA we can use:
$ kubectl get hpa

Thank you.