vol types

1) Node local type such as ~~empdir~~ emptydir
   empty dir and hostpath.

ii) File sharing such as NFS

iii) Cloud provider, Aws eBS or Azure disk

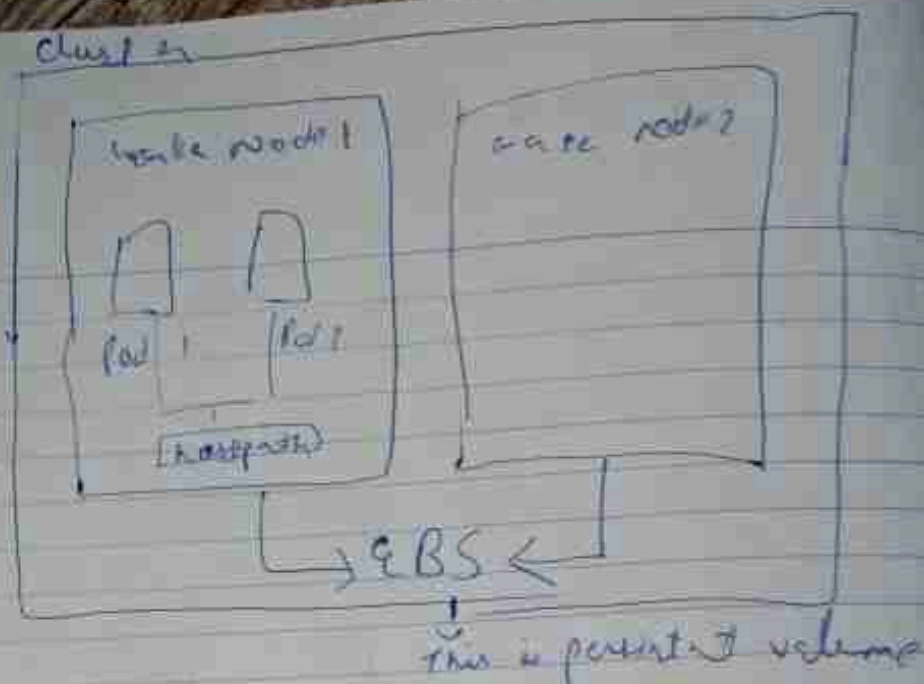iv) distributed file system glusterfs or cephfs

v) special type secret, git repo


* PERSISTENT VOLUME  AND  LIVENESS PROBE

A persistent volume (PV) is a cluster - wide
resource that you can use to store data
in a way that it persists beyond the
lifetime of a pod.

The PV is not backed by locally attached storage
on a worker node but by networked
storage system such as EBS or NFS or a
distributed file system like ceph.

K8s provides APIs for users and administrator
to manage and consume storage. To manage
the volume, it uses the persistent volume
API resource type and to consume it, uses the
persistent volume - claim API resource type

cluster

worker node 1 — worker node 2

Pod 1   Pod 2

[hostpath]

→ EBS ←

This is persistent volume

# PV Claim

- The PVC request a PV with your desired specification (size, access type) from the and once PV is found, it is bound to PV claim.

- After a successful bound to a pod, you can mount it as a volume.

- Once a user finishing it works, the attached PV volume can be released. The underlying PV can there be reclaimed and recycled for future usage.

## Restrictions using AWS EBS

1) The nodes on which pods are running must be AWS EC2 instance.

2) The instances need to be in same region and Availability zone as EBS volume

3) EBS only support a single EC2 instance mounting a volume

# HEALTH Check / liveness probe

- Health checks or probes are carried out by the kubelet to determine when to recreate a container (for liveness probe) and used by services and deployment to determine if a pod should receive traffic.

- One use of readiness probe is to control which pod are used as backends for service. When a pod is not ready, it is removed from service load balancer.

initial delay second = time do container that have ke lige
period seconds = agla health check kab karna hai
timeout seconds = kill karo ehorwing karna aur container restart karo

# Configmap AND SECRETS

Configmap can be created in following way :-

1) As env variable
1) As vol in the pod

Kubectl create configmap <mapname> --from-file
                = <file to read>

# SECRETS

1) Secrets are namespaced objects, that is exist in the context of a namespace.

7) The secret data on nodes is saved in tmpfs volumes. tmpfs is a file system which keeps all file in virtual memory. Everything in tmpfs is temporary in sense that no file will be created on your hard disk.

8) You can access via env variable or volume

4) A per secret size limit of 1 MB exists

5) The API server stores secret as plaintext in etcd

* Secrets can be created

i) from a text file
ii) from a yaml file.

* to make a secret,

• kubectl create secret generic mysecret --from-file = volume.txt --from-file = pwd.txt

• kubectl get secret.

1 CPU = 1000 m

memory = miB

x NAMESPACE , Limits ad Requests

• A namespace is a group of related element
that each have a unique name or identifier
Namespace is used to uniquely identify
one or more names from other
in to names of different objects, groups or
the namespace in general

• A scope for every names

• A mechanism to attach authorization and policy to a
subsection of the cluster.

• Most K8s resources ( pods, services, replication
controller) are in some namespace ad
low level resources such as nodes and
persistent volumes are not in any namespace

* To change from default namespace,
kubectl config set-context & (kubectl config
current-context ) -- namespace = dev

* To create a pod in your namespace
kubectl apply -f pod.yml -n dev ← namespace
name

✦ Managing Compute Resources :—

• Scheduler decides about which nodes to place
pods , only if the node has enough
CPU resources available to satisfy the
pod CPU required

* Two types of Constraints can be set of each resource type

1) Request and 2) Labels

1) A request is the amount of that resource that the system will guarantee for the container and k8s will use this value to decide on which node to place the pod.

2) A limit is the max amount of resource that k8s will allow the container to use. In the case that limit is not set for a container, it defaults to none, if limit is not set, then if default to 0. (unlimited)

| | |
|---|---|
| Request = not nothing | Request = nothing |
| Limit = nothing nothing | Limit = no nothing |
| Request = limit | limit = default |

You can limit
1] Compute, memory, Storage.

2 Restrictions

1) Every container that runs in namespace must have its own CPU limit

2) Total amount of CPU used by all container in the namespace must not exceed a specified limit.
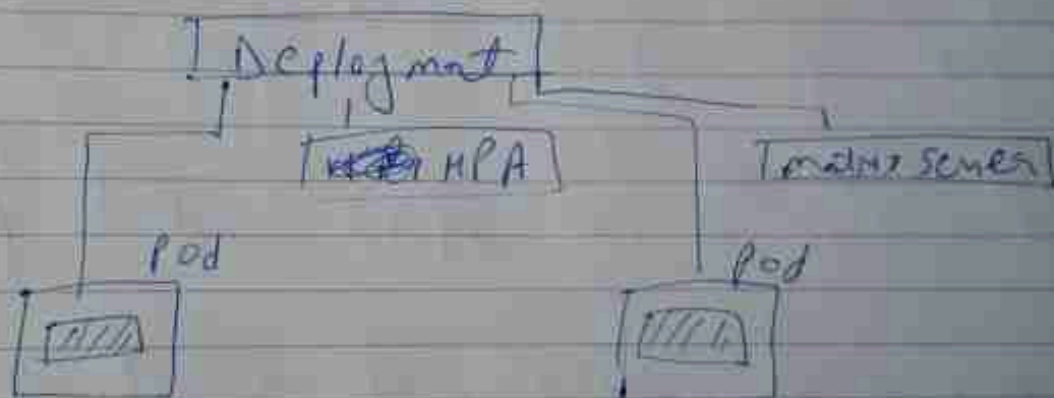
Resource Quota and Horizontal Scaling

Default Range →

CPU
(min) request - 0.5
(max) limit - 1

Memory
(min) request - 500 M
(max) limit - 1 G

• Horizontal Pod Autoscaler

• KBs has the possibility to automatically scale pods based on observed CPU



cooling period = 5 min
HPA check in every 30 second

○ Scaling can be done only for scalable objects like controller, deployment or replica set

• HPA is implemented as a KBs API resource and a controller

• HPA is implemented as a controlled loop with a period controlled by the controller manager.

horizontal pod autoscaler loop period

## Jobs

Jobs gets stop after doing his work.

Objects like replicasets, daemon sets, statefulset and deployments they all make sure their pod are running

Cron Job
If we have multiple nodes hosting the application for high availability, which nodes handle cron?

* Init Container.
  - Init Container are specialized Container that run before app container in a pod.
  - Init containers always run to completion.
  - Init container do not support readiness probe.

  use case,
  1) Clone a git repository into a volume.
  2) Generate configuration file dynamically.

* POD LIFECYCLE

| Pending | Running | Succeded | failed | Completed | Unknown |
|---------|---------|----------|--------|-----------|---------|

1) Pending.

- The pod has been accepted by the K8s system but it's not running.

- If the pod cannot be scheduled because of resource constraints

2) Running
- The pod has been bound to a node.
- All containers have been created

3) Succeed
All containers in the pod have terminated
in success and will not be restarted.

4) Failed
All containers in the pod terminated and atleast one
container has terminated in failure.

the container either exited with non-zero status
as was terminated by system.

5) Unknown
- state of the pod could not be obtained

- Typically due to error in retrieval or communicating
  with the host of the pod

6) Completed
the pod has run to completion as there's
nothing to keep it running

◊ POD conditions

Kubectl describe pod <Pod name>

1) Pod Scheduled : Pod has been scheduled
2) Ready : added to load balancing pool
3) Initialized pod : All init containers have started successfully
4) Unscheduled : can't schedule right now
5) Container Ready : All containers are Ready