

## ● Access nginx-svc by using Name :-

Name of my Service: nginx-svc

```
# curl http://nginx-svc
```

o/p. <!DOCTYPE html>

<html>

<h1> welcome to nginx! </h1>

</html>

So, its working.

# nslookup nginx-svc → nslookup Command

↓  
It gave IP address: 10.0.2.134. what it does?

But other things it didn't find.

Q) what does nslookup do and what is purpose of /etc/resolve.conf

Ans: Nslookup (Name Server lookup) is a command for getting info. from the DNS Server.

It is a network admin's tool for querying the Domain Name System (DNS) to obtain domain name (a) IP address mapping (b) any other specific DNS record.

It might be used to figure out what your name server is?

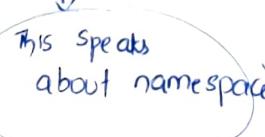
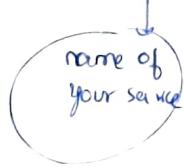
If you give name — it will give you IP-address.

→ There is a concept called as Fully Qualified Domain Name → is a complete name of your service.

```
# nslookup nginx-svc
```

Full name: → nginx-svc.default.svc.cluster.local.

Address: 10.0.2.134.



↳ is a default domain name inside your k8s cluster.

If you want → You can change that.

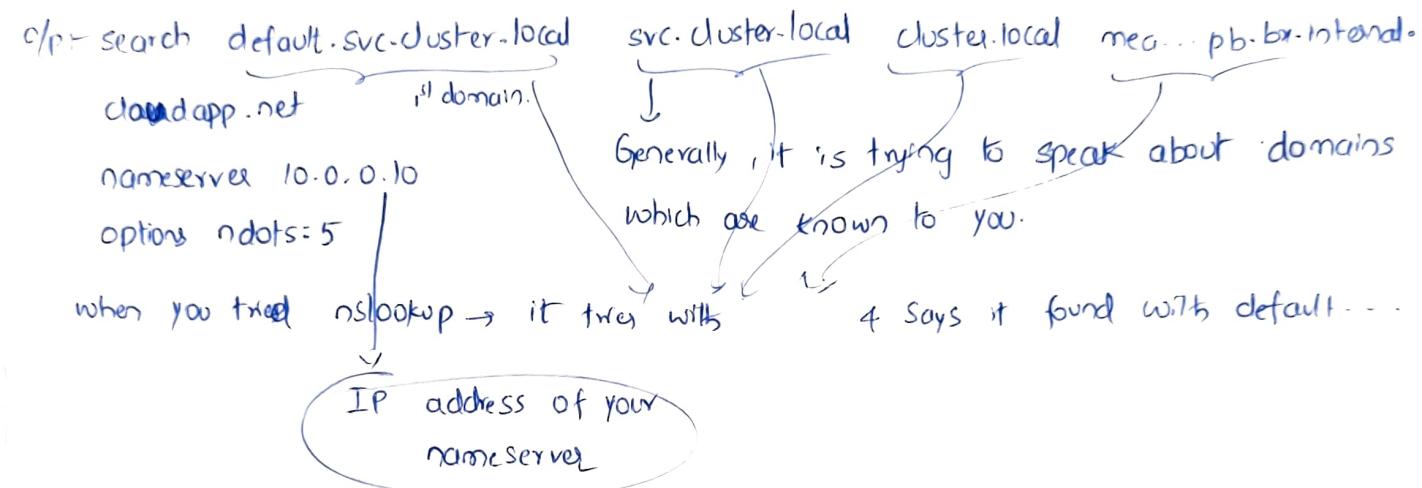
/# nslookup 10.0.242.134

O/p:- It gives name of your Service : nginx-svc.default.svc.cluster.local.  
So, are there DNS entries.

So, you can access services by using its name (or) IP-address. Both are not going to change.

/# cat /etc/resolv.conf

This consists of local info. about your domain names and name servers.



Look into Environment Variables in alpine pod ( Alpine was created post nginx service creation )

/# printenv

O/p:- KUBERNETES\_SERVICE\_PORT=443

1<sup>st</sup> we created → nginx-rs  
2<sup>nd</sup> → nginx-svc  
3<sup>rd</sup> → alpine-pod.

HOSTNAME=alpine

HOME=/root

NGINX-SVC-SERVICE\_PORT=80

→ this is the name of our service.

NGINX-SVC-PORT=tcp://10.0.242.134:80

\* So, k8s will inject Environment Variables about your service into every Pod (container) that gets created after container creation.

That means

In k8s, after we create a Service → if we create a container → it will have Environment variables related to service. But these are not dynamic. So for eg. if you remove the service these variables will not be removed. So relying on them is not a good idea, but k8s does this

\* Look into Environment Variables in engine pods (They were created prior to nginx-svc) :-

(Linux VM → which host)

Let's use other machine kubectl

(53)

\$ kubectl get po

% alpine

nginx-rs-6gjfb

nginx-ss-zqxrld

\$ kubectl exec ~~①~~ nginx-rs-6gjfb -- ~~②~~ printenv

O/p:- HOSTNAME=nginx-rs-6gjfb

KUBERNETES\_SERVICE\_PORT=443

Relying on Env. variables is not a good idea. (It doesn't give dynamic info)

Relying on nslookup is a good idea (it gives dynamic info)

HOME=/root.

In this, are you seeing nginx-service? → No.

Bcz, this ~~①~~ was created before service's creation.

/H Exit.

Service  
means ?

\* Cluster IP → gets an Internal IP address within K8s cluster. It can't be accessed to the outside world.

Google: Kubernetes Services. ( Publishing Services : Service Types )

\* Now, what I want to do is → I'm not interested in alpine pod.

So, I have some user who is outside → who wants to access nginx-service.

This person who is outside → how can he access nginx-service?

## External Communication using Kubernetes Service :-

\* Some user external to K8s cluster wants to access nginx.

\* For this what K8s tries to tell is → I have different types (Service Type)

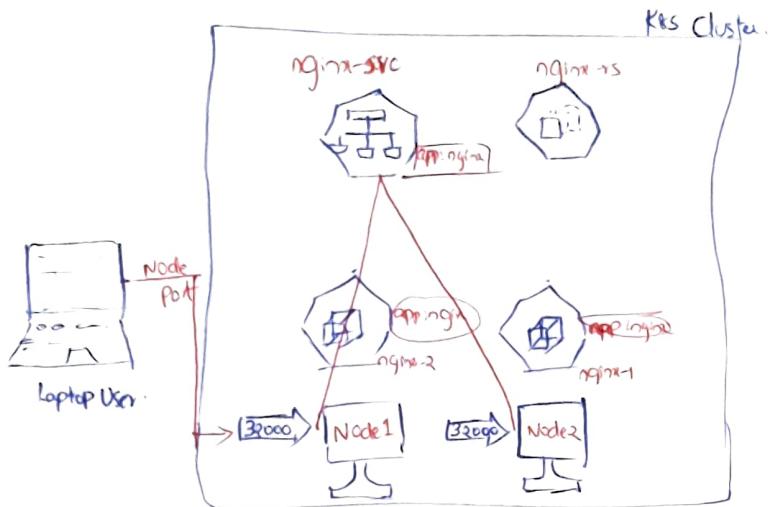
Right now, the only type which we know is → ClusterIP → which is an Internal IP address

→ K8s has the service publishing types :-

1) ClusterIP - <sup>for</sup> Internal Comm?

2) Node Port → K8s will expose the application on a port on every node in K8s cluster

NOW, when you access 32000 port → This basically connects to nginx-service



This is like Container  
Port-forwarding

But, in Q1 we had only one M/C.  
Here we have multiple M/C's.

So, what we are trying to tell is

Your system needs to basically have a connection to one of the nodes and it  
wud try to access the nodes on some port → and that wud be accessing the app?

For this to work → your nodes need to have public-ip address. , if they want to  
be accessed from Internet.

(You don't want to expose the IP-addr of your node)

(This is useful within On-premises).

Till yesterday we were using kubeadm → Remember all three nodes had public ips)

\*Now

If you want to expose nginx-SVC, you will be using node port.

Execute

Try on kubeadm

## \* Load Balancer :-

This is generally used with managed K8s clusters .

Here, what we are trying to tell is. →

Every cloud has some kind of Load Balancer (Layer4 Load Balance)

Whenever you create service type as Load balancer →

In case of Azure → it creates a Azure load bal.

This job is done by Cloud-Controller Manager (CCM) (bcz K8s by default will not  
have intelligence of azure)

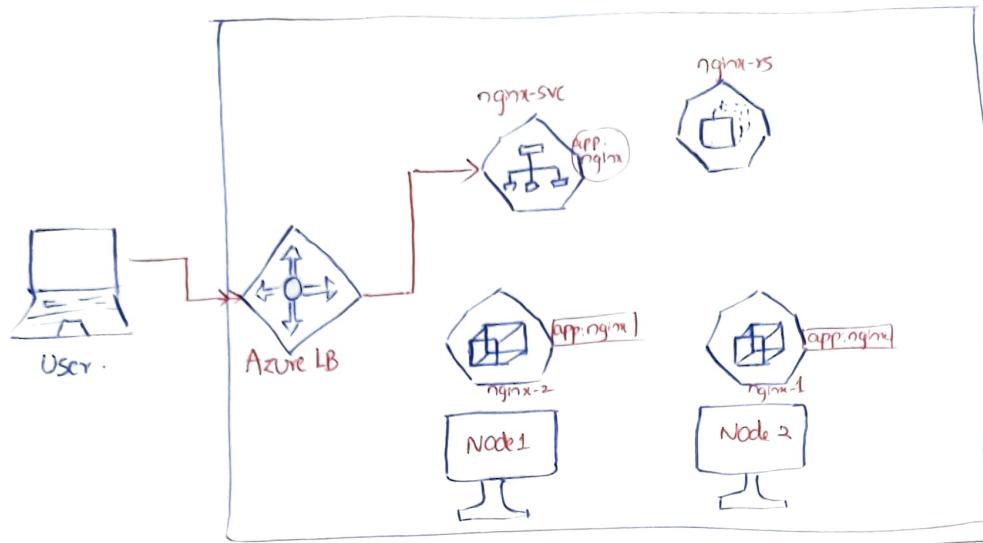
Now You will be accessing load-balancer → which will forward the request to nginx-SVC.

This is the 3rd type.

also a load balancer

This is decent way. But there is a better way.

(54)



ClusterIP is  
Private IP address

The type of service is : Loadbalancer

→ This will not work in kubeadm as it doesn't have CCM. It doesn't know how to create a LB. If you are using kubeadm, the only thing which you can resort to is 'node port'

4th type: External :-

Service-type

Here you get a CName record



CNAME record is something which you can put in DNS servers.

2nd type

When you basically put an external name → you get a cname → Add the entry to your DNS server which is maintained by your orgn so that commun' goes through.  
(we use this rarely).

\* Node Port and External name are used in on-premise k8s clusters a lot.

In the case of cloud → basically we tend to use Load balancers.

- † Service types :
- 1) ClusterIP
  - 2) nodeport
  - 3) Load balancer
  - 4) ExternalName.

'SVC → mkdir .\lb-demo. → nginx-lb-svc.yaml.

\svc\lb-demo:>

---

apiVersion: v1

Kind: Service

metadata:

name: nginx-lb

Spec:

ports:

- name: nginx-lb
- port: 80
- protocol: TCP
- targetPort: 80

selector:

app: nginx

type: LoadBalancer.

Diagram annotations:

- Ports: points to ports section
- Service Ports: points to targetPort
- name: points to name
- port: points to port
- targetPort: points to targetPort
- Loadbalancer's Port: points to targetPort
- targetPort → is where your app is running in the @'s

> kubectl apply -f .\nginx-lb-svc.yaml.

> kubectl get svc -o wide -w

IP	Name	Type	cluster-IP	External-IP	Port(s)	Selector
	Kubernetes	clusterIP	10.0.0.1	<none>	443/TCP	<none>
	nginx-lb	LoadBalancer	10.0.118.52	20.121.176.189	80:30087/TCP	app=nginx
	nginx-svc	ClusterIP	10.0.247.134	<none>	80/TCP	app=nginx

We got External IP over here.

It is actually IP-address of your LoadBalancer.  
(Go to Azure → Load Balancer → )

Let's access using this IP-address

http://20.121.176.189

O/P: ! Welcome to nginx!

The port on which you'll be accessing app is 80.

If I have given 8080 → then it would be

20.121.176.189:8080

which will map to 80 port of nginx.

\* So, whatever applications you have written → let it be e.g. nopCommerce, game of life, Spc etc

Can we expose it to outside world now? → Yes.

All you need to do is

→ If you are using kubeadm → use nodePort.

If you are using Azure AKS → use type as LoadBalancer.

This is not a better way.

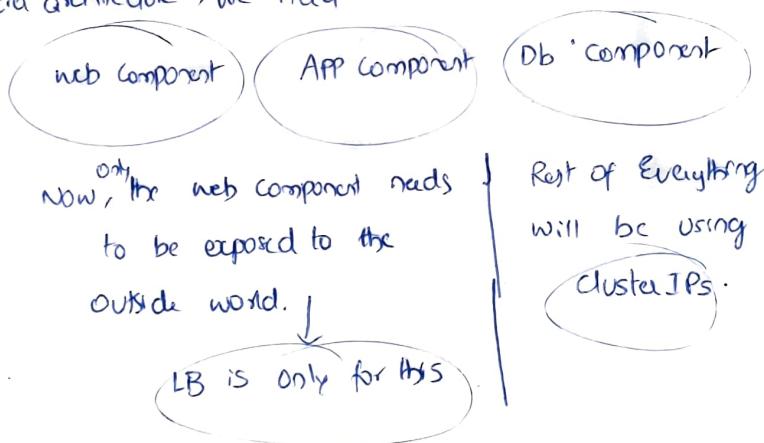
There is something more which we are going to learn.

In that we will try to access the app in a much better way.

For now, this is ~~is~~ a Layer-4 LB.

(Layer 7 is better than Layer 4).

\* In n-tier architecture, we had



\* Layer 4 & Layer 7 are OSI layers.

→ I don't want to put my Pod in running state until its dependency is up.

① (I have app server & I have db server → I don't want to put my app pod in running state till db is up).

② In Pods, I'm running my apps. Apps might fail. → So when it fails, what should we be doing? How should we be configuring our checks?

Deleting the Cluster → az group delete --name myResourceGroup --yes --no-wait.

# L-56-Apr-29-Part-2 | Health Checks/Probes for Pods in Pods, Container Types in Pods

(56)

2 tabs

1) PowerShell → Laptop    2)    Linux kubectl.

\$ ssh ~~ubuntu~~ 52 25.224.255. (Azure VM)

at kubectl

→ azure aks quickstart (lets quickly create k8s cluster).

\$ az group create --name myResourceGroup --location eastus.

\$ az aks create -g myResourceGroup -n myAKScluster --enable-managed-identity  
--node-count 2 --enable-addons monitoring --enable-msi-auth-for-monitoring  
--generate-ssh-keys.

We have created nginx-rs & nginx-pods.  
 ↓  
 4 replicas.

Label → app:nginx  
version → v1.23

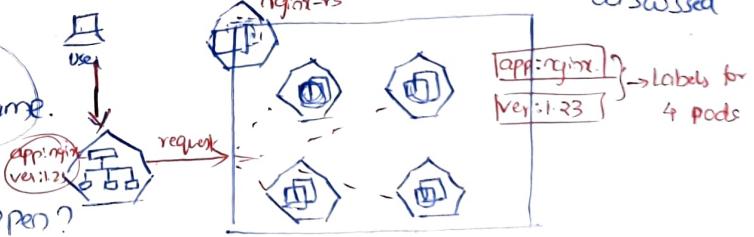
I have created a service. (Label: app:nginx)  
now version: 1.23

\* whoever wants to make a request. (Someone is raising the request).

\* k8s service will forward the request to 4 of them. (This is what we discussed)

Now what has happened is?

Our application's new version came.



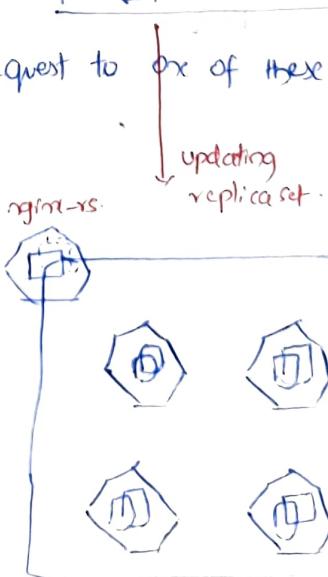
\* So right now - what would happen?

Your nginx-svc will try to forward the request to one of these 4.

Then, what we have done is



we have updated our replicaset



\* Initially what would happen is

it will forward the request to upper (1.23)

Now, when you update and you have a label  $\xrightarrow{\text{nginx}}_{1.24}$ , all of these pods (1.23) will be dead.

Now if someone accesses k8s Service  $\rightarrow$  it will not be able to forward the request bcoz  $\rightarrow$  when you update (all the pods are gone). You don't stop older pods cos so, 1.23 version  $\rightarrow$  there is no pod.

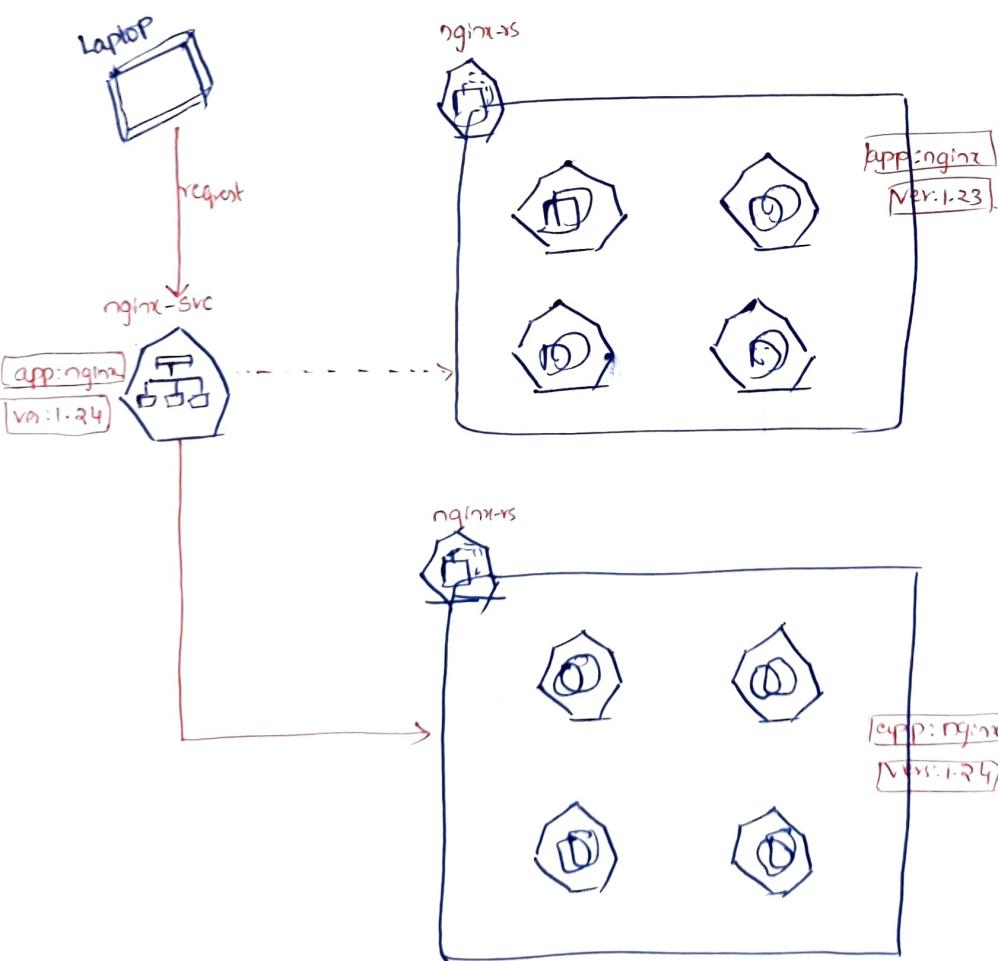
↓  
replica set.  
You create one more rs.

Add label ver:1.24 to a k8s-Service. (then previous link will be broken).

\* Generally, whenever you want to update to a new version, without having a downtime, this is the approach we follow:

\* when upgrading to newer versions of Pods, ensure right set of labels are present on k8s Service ~~else~~ Selector.

You will set the label  $\rightarrow$  such a way that it will forward the request to older ones as well as newer ones. One if you think everything is working  $\rightarrow$  replace it with version in SVC.



\* But, the problem is, if your application is running inside k8s → that means we are running app<sup>n</sup> inside pod → indirectly, you are running app<sup>n</sup> inside the Container.

But is it guaranteed that every app<sup>n</sup> inside the Ⓛ will work without any problem? → NO.

• So, when your app<sup>n</sup> is not working → k8s shud know that. So, how will k8s know the component → what is the k8s component which reports the status of something running on the worker node to master node?

↓  
Component is → Kubelet → ~~pod~~ shud know what is happening in your pod.  
for that, let us try to understand

↓  
probes for Ⓛ's  
\* Health Checks/ in k8s Pods :-

Configure kubelet to connect to your k8s cluster using az aks get-credentials command.

\$ az aks get-credentials --resource-group myResourceGroup --name myAKScluster.

y  
y  
Dell> az aks get-credentials ..  
y  
y  
y

Google: health checks in k8s.

• So k8s supports 3 kinds of checks

- 1) Liveness Probe
- 2) Readyness Probe
- 3) Startup Probe

For every k8s Ⓛ → we can try to write all 3.

The basic idea is → liveness is used for the k8s to know when it shud restart the Container.  
we remember that when k8s goes into exited state → k8s will restart the Ⓛ.

• Liveness Probe :- Liveness is used for the k8s to know when it should restart the  $\textcircled{O}$ .

But, let us assume that, your app has become too much busy, some process has become still  $\rightarrow$  restarting the  $\textcircled{O}$  might solve the problem  $\rightarrow$  So for that, if you try to write a liveness probe (liveness check) and then if that check fails  $\rightarrow$  we want k8s to restart the  $\textcircled{O}$ .

That is exactly what liveness Probe is all about.

• Readyness Probe :- Readyness Probe  $\rightarrow$  what it does is  $\rightarrow$  if Readyness Probe is passing, then only your service will forward the request to the Pod.

If Readyness check fails, your  $\textcircled{O}$  might be up & your pod might be up  $\rightarrow$  but it will still not forward the request.

K8s Service will not consider - your pod - to forward the request until & unless that check becomes passed.

Kubernetes has recently introduced  $\rightarrow$

• Startup Probe :- is just to figure out whether everything in your application is started (or) not.

Till Startup Probe is finished, it will not check for liveness & readiness.

Liveness Probe: If this check fails, K8s will restart the  $\textcircled{O}$ .

Readiness Probe: If this check fails, the Pod will removed from Service (Pod will not get requests from Service).

Startup Probe: This checks for startup and until startup is OK, the other checks will be paused.

It is always a good idea to write Health checks for your  $\textcircled{O}$ 's  $\rightarrow$  b/cos you can tell when your App is ready, when your App has started and when you think that you shud restart your App.

• To do these Health checks,

- API reference  $\rightarrow$  Pod v1 Core  $\rightarrow$  PodSpec  $\rightarrow$  Containers  $\rightarrow$  Probes -  
 $\begin{array}{l} \text{exec} \\ \text{grpc} \\ \text{httpGet} \\ \text{tcpSocket} \end{array}$

\* Probes (or) checks can be performed by

- exec: run any linux/windows command which returns status/exit code.

We all know that if exit code is zero(0) → it is success.

If it is anything other than zero → it is failure.

So, to know the health of your app — if you have to run some command — that is what you do using exec.

- http: We send http request to the app.

Any http requests will give you status codes (3 digit no.)

Based on status codes we can decide.

Google: http status codes.

- Grpc: Is also one of the communication mechanism

This communicates over grpc.

starts with 1 → i.e. Information

" " 2 → Success

" " 3 → Redirection

4 → Client-side Error

5 → Server-side Error.

In k8s, if you are trying to run a pod which is a micro-service developed work — which exposes grpc — then you would be using this grpc to figure out whether your app is healthy or not.

- tcp: send tcp request.

~~http~~ folder → healthchecks folder → httpd.yaml

To cross-check

Let us try to create any simple app →

F12 → Network

when you type → fb.com

it goes to facebook.com.

Status code: 301 (redirection)

There is nothing called as fb.com.

E1: directdevops.blog → status 200

E2: directdevops.blog/testing123 → success

Error → 404.

(client-side error)

- Now, we will be writing multiple yaml files in one file.

Here we will create 2 specs

Service

replica set.

# --- httpd.yaml (Service file)

apiVersion

    v1

metadata

    name → httpd-svc

spec

    ports (Service Port)

        ↓  
        port ( port that will be exposed by this service → 80.)

            name → httpd-svc.

            target port → 80.

selector

    app: httpd.  
    type → LoadBalancer.

You can use clusterIP & use curl requests  
for checking

replica set

    ↳ 2 replicas of httpd:-

version → apps/v1

kind → Replicaset

metadata

    name: httpd-rs

    labels:

        app: httpd.

spec:

    ↳ replicaset Spec.

        ↓  
        minReadySeconds: 1

        replicas : 2.

        selector:

            matchLabels:  
                app: httpd.

        template:

            metadata .

                name → httpd

                labels:

                    app: httpd.

spec :

    containers:

        - name: httpd

        image: httpd:2.4.57

    ports:

        - containerPort: 80 .

    livenessProbe

        Type: (Probe)

            click on it

        exec (or)

        grpc

        failure threshold  
        = 3 .

        Initial DelaySeconds

        = 1

Default .

After Success .

(It will wait for 3 consecutive failures)

This check is considered to be failed  
when it fails 3 times.

we can change that no → 1 (can't put 0)

When do you want to start probes

PeriodSeconds  
= 10

default

do you want to start

Immediately once the container has gone  
into running state

How frequently do you  
want to do the probes

(or)  
after some time )

SuccessThreshold = 1

↓  
how many checks should it pass to consider it a Success.

Default = 1.

timeoutSeconds = 1

↓  
Once you make this check → by what time do you expect your check to get response.  
(you have sent a request)      "      "      "      response to come

httpGet

(HTTPGetAction)

↓ click on it.

Default = 1

where I should get  
200 as code

There is no need for you to give host  
as you are taking → local host.

Path = / → means  
homepage.

Port: 80. → on which our app is running

For some reason if this app goes for a  
loss → now K8s will try to restart.

readiness Probe:

(Probe).

write same thing you have written in live probe  
as both are of same type.

- We will write failure scenario

↓ writing  
exec:  
command  
- sleep  
- 10 MM

→ in livenessprobe  
instead of error → failed scenario.

HttpGet:  
Path: /  
port: 80

Exec action  
2nd method

3rd one is

→ sending TCP requests

4th " → sending gRPC requests.

Lets try to run this

docker playground.

↓ new instance

→ docker container run -d -P

--name apache httpd:2.4.57

docker container ls

O/P: ... httpd:2.4.57 ... 0.0.0.0:49153 > 80/tcp

\$ curl http://localhost:49153

9P <It works>

We are not interested in response.

We are interested in status code.

Google: curl <sup>linux cmd</sup> examples.

Curl to print status code.

\$ curl -I http://localhost:49153

9P: HTTP/1.1 200 OK. ✓

O/P

\$ docker container exec apache

and try to run some Linux cmd which  
will check whether apache process is  
running or not.

If apache process is running → it will  
give return code 0 (otherwise 1).

> cd Manifests\rs\healthchecks

> kubectl apply -f ⑥

O/p: Pod/alpine created

Service/httpd-svc created

replicaset.apps/https-rs created

If you give folder name—it will try to run everything in that folder.  
↓  
we have 2 files.

> kubectl get po

	<u>status</u>	<u>Restarts</u>	<u>Age</u>
alpine	1/1 running	0	69s

http-- 1/1 "

2

"

→ here we gave

timeout seconds=10

http-- 1/1 "

2

"

we are getting restarts (due to failure)

change to 1.

> kubectl apply -f .

O/p: nothing changed.

Only rs-configured.

> kubectl get po -w

alpine	1/1 running	0
http--	1/1 "	3
"	1/1 "	3

> kubectl describe pods https--

O/p: Logs → will speak about Liveness  
Events. Readiness

↓ will show → Liveness probe has failed.

so I'm supposed to restart my app? (Actually my O/P has not gone)

into excited state but  
still I need to restart  
(blog liveness check-failed).

Now replace with

httpGet:  
path: / in liveness.probe.  
port: 80

> kubectl apply -f .

> kubectl get po -w

O/p: from now restarts should not happen

bcz our checks will pass now.

httpd.yaml :-

```
apiVersion: v1
kind: Service
metadata:
  name: httpd-svc
  labels:
    app: httpd
```

```
spec:
  ports:
    - name: httpd-svc
      port: 80
      targetPort: 80
  selector:
    app: httpd
  type: LoadBalancer
```

```
apiVersion: apps/v1
```

```
kind: ReplicaSet
```

```
metadata:
  name: https-rs
  labels:
    app: httpd
```

```
spec:
  minReadySeconds: 1
  replicas: 2
```

```
  selector:
    matchLabels:
      app: httpd
```

```
  template:
    metadata:
      name: httpd
      labels:
        app: httpd
```

## Spec :

## containers :

- name: httpd
 image: httpd:2.4.57

## ports :

- containerPort: 80

## livenessProbe :

- failureThreshold: 3

- initialDelaySeconds: 1

- periodSeconds: 10

- successThreshold: 1

- timeoutSeconds: 1

## httpGet :

- path: /

- port: 80

## readinessProbe :

- failureThreshold: 3

- initialDelaySeconds: 1

- periodSeconds: 10

- successThreshold: 1

- timeoutSeconds: 1

## httpGet:

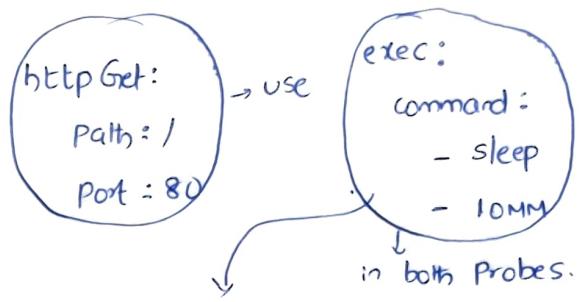
- path: /

- port: 80.

## http-faulty.yaml :-

Same as previous file except → instead of

Replica-set name: https-rs-faulty.



what would happen?

Here we are running 2 pods → (we already have 2 healthy pods → httpd.yaml).

Generally when you write this,

service should actually serve the 4 pods. It should forward the request to 4 pods.

### Readiness

We have said in Service → selector → app: httpd.

So, here also in pods we have app: httpd (httpd.yaml)

whatever is faulty → there also i have app: httpd (http-faulty.yaml)

→ But, the point is, then pods should not be served bcoz Readiness is failing

→ Ideally, your service should forward the request to 4 pods. But it will be serving to

Only two.

> kubectl apply -f .

O/p: Required value → must specify a handle type. (Actually we are not worried about liveness probe.

we want service to only serve 2 pods by failing → 2 pods in faulty.yaml

	cluster-ip	External-IP
o/p: httpd=svc	-	20.88.179.76
Kubernetes	-	-

> kubectl describe svc httpd-svc.

O/p: - -

Endpoints: 10.244.0.10:80 , 10.244.1.12:80.

- -  
we are seeing Endpoints.

Endpoints are your (O) IP-addresses → your service is serving.

When we create Service, end points are created

↓  
End pt represents each (O) pod's IP-address.

\* So, totally we have 4 Pods which have the matching labels  
 But totally, how many is it serving? = Q. ( busy in other ? (faulty) → Readiness failed)

\* Now to correct the `http-rs-faulty.yaml` → replace

> kubectl apply -f .

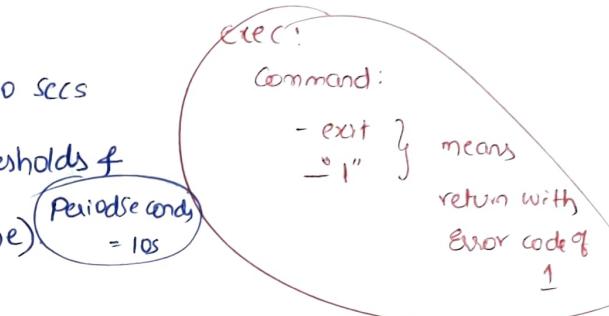
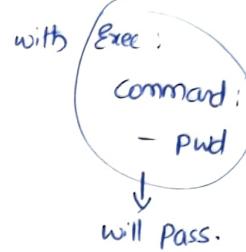
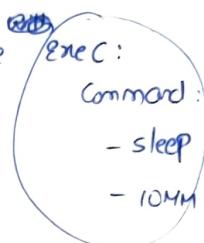
> kubectl describe svc httpd-svc

> kubectl get endpoints

After changing exec and apply → wait for 30 secs

because → In Probes → we have written success thresholds of failure thresholds. (wait for checks to be done)

(wait for the restarts)



For a check to fail

→ it shud fail 3 times after a success.

Threshold for failures is 3 after success

$$3 \text{ checks} = 3 \times 10s = 30 \text{ sec.}$$

(Liveness probe is all about issues that may be resolved by using restarts)

\* From now on, whenever we write Pods,  
 we will try to write some necessary probes.

## Run Pods with Specific Resources (CPU/Memory) :-

So, How much of CPU & Memory are you speaking about?

→ I want to run  $\odot^s$  in my pod with specific resources → I want to give them what is the minimum CPU which it shud consume & what is the maximum? what is the minimum RAM it shud consume & also Max

busy → lets assume that → ~~on historical basis~~ On a historical basis, our Orgn is basically running a NYC with 4 CPUs & 16 GB of ram and there we can run 100 pods.

But some app is misbehaving and using too much of CPU. so → just bcz ~~that app~~ is behaving → and too much of CPU → You are impacting other app's bcz CPU of that app goes for a toss.

How to limit that?

Exercise

# Try running a docker Container with 256mb of ram & 1 CPU.

Run any Docker → Open Docker playground & figure out a command for that.

Type the command → docker stats.

Ans. --cpus

--memory.

Eg: docker run -d -p 8080:80 --memory="256m" nginx

docker run -it --cpus=".5" ubuntu /bin/bash

If you have 1 CPU, the command guarantees the Docker at most 50% of the CPU every second.

(Or)

docker run -it --cpu-period=100000 --cpu-quota=50000 ubuntu /bin/bash

Google: Resource Management for Pods  
Manage Resources - containers

Here, in K8S, whenever you want to put restrictions for resources' usage → there are

2 fields → 1) Requests → is the minimum amount which you want.

2) Limits → is the maximum amount which your pod can go.

• for Eg: If I want to run a Pod → min of 256mb, max of 1Gb.

whenever you are running a Docker in a Pod.  
It is always a good idea to put restriction

I want to give  
will go in Requests  
will go in Limits.

You can also skip one. (if you can write only limits)

• These will be part of Container Spec.

→ rs folder → limits folder → nginx-rs.yaml :-

---

apiVersion: apps/v1

kind: Replicaset

metadata:

name: nginx-rs

Spec:

minReadySeconds: 1

replicas: 3

selector:

matchLabels:  
app: nginx

selector

template:

metadata:

- : name: nginx-pod
- : labels:
- : app: nginx

Spec:

containers:

- name: nginx
- : image: nginx:1.23
- : ports:
- : - containerPort: 80

livenessProbe:

- : httpGet:
- : path: /
- : port: 80

readinessProbe:

- : httpGet:
- : path: /
- : port: 80

resources:

- : requests:
- : memory: "64Mi"
- : cpu: "250m"

limits:

- : memory: "256Mi"
- : cpu: "1000m"

apiVersion: v1

kind: Service

metadata:

- : name: nginx-svc

spec:

selector:

- : app: nginx

type: ClusterIP

ports:

- name: nginx-svc
- : port: 80
- : targetPort: 80
- : protocol: TCP

> kubectl delete -f .

Writing a Probe & }  
Writing a Resource limits } is always a good Practice.

> cd limits

> kubectl apply -f .

> kubectl get po

> kubectl describe pods nginx-ss-k8bnt

O/P: Limits:

CPU : 1

Memory : 256 Mi

Requests :

CPU : 250 m

Memory : 64 Mi

## Container Types in Pods :-

- In Pods, we have 3 types of Containers

- 1) Containers: These are why we write Pod spec
- 2) init Containers
- 3) ephemeral containers.

Google: init container kubernetes.

> kubectl delete -f .

→ Init Containers :-

When K8S starts Pods → In Pods you are starting ○'s

- First it will start init ○'s.
- Init ○'s are the ○'s which should not run forever. They should be completed if you remember Job → Job is a ○' which basically finishes in some time.  
So init ○'s are the ○'s where you do Pre-condition.

1000 m ≈ 1 CPU

250 m =  $\frac{1}{4}$  of 1 CPU

→ For Eg: for running your application, if you basically need certain dependency - other  $\circ$ s shud be up or not.

Till other  $\circ$  is up → don't start the  $\circ$ . If this is your use-case - then what we will try to do is →

we will try to put a init container in which we will wait for the other service to be up.

↓  
Only after the other service is up → then we will finish the init  $\circ$ .

Once we finish the init  $\circ$  → then <sup>your</sup> actual  $\circ$ 's will start in which your app runs.

- So, the basic flow is :

Pods will create init  $\circ$ 's.

In init  $\circ$ 's you can write 100 containers if you want → but it will run one after other. It will not try to start them in parallel.

(In pod → if we have created 2  $\circ$ 's → it will start 2  $\circ$ 's at the same time)

(But init  $\circ$ 's are not like that).

- Till init  $\circ$ 's have reached its success state → Your actual  $\circ$ 's will not be started.

\* If your app has some pre-conditions to be met

↓  
so, 1<sup>st</sup>, we'll be writing those pre-conditions as init  $\circ$ 's.

\$ dig google.com

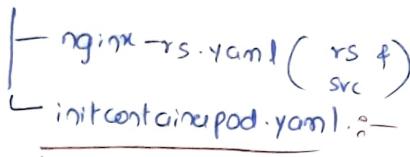
This dig command will basically try to run and it will try to tell you what is a query time and some stuff.

But if the service is not running.

↓  
\$ dig google1.com → Here, it will basically try to give you Errorstate

Eg: for i in {1..100}; do sleep 1s; if dig myservice; then exit 0; fi; done;

DS folder → Init containers



apiVersion: v1

Kind: Pod

metadata:

name: init-demo

Spec:

containers: → my actual container over here.

- name: httpd

image: httpd

ports:

- containerPort: 80

InitContainers:

- name: delay

image: alpine

command:

- sleep

- 15s

- name: init-myService

image: busybox:1.28

command:

- 'sh'

- '-c'

- "until nslookup nginx-svc.default.svc.cluster.local; do echo waiting for myservice;"

↓

I'm trying to run a script. → what this script will do is → it will try to wait till this service is up. and only when then source is up → then your actual containers will be created.

Pod

PodSpec

↓

→ Containers

→ ephemeral C

→ init Containers

(type: Container Array).

docker playground

↓

> ssh ...

\$ docker container inspect apache

O/p: - - -

IP address 172.17.0.2

\$ docker container run -it alpine /bin/sh

# for i in {1..100}; do sleep 1;

if dig 172.17.0.2; then exit 0; fi;  
done; exit 1.

O/p: dig not found.

# apk add dig bind-tools.

# for i in ... = It is failing

if it exits out of this → it's good.

\$ docker container run -it alpine/bin/sh

# for i in {1..100}; do s... 172.17.0.10.

Let me give  
something which  
doesn't exist →

sleep 2; done"

1st it will sleep for 15 seconds and then it will wait for other service to be up  
and then

once that is done → our actual container will start.

> cd initContainers

> kubectl apply -f .

O/p: pod/init-demo created  
replicaset.apps/nginx-rs created  
service/nginx-svc created.

> kubectl get po -w

	<u>Ready</u>	<u>Status</u>	<u>Restarts</u>	<u>Age</u>	
init-demo	0/1	Init:0/2	0	9s	In Init it has to finish 2 containers (15 sec) → (1st). nslookup.
nginx-rs-ct.	1/1	0	0	8s	
nginx-rs-gf-	1/1	0	0	9s	
nginx-rs-wt	1/1	0	0	8s	It has not gone into running state immediately.
init-demo	0/1	Init:Y2	0	16s	It is waiting for 2 init O's to finish its execution
init-demo	0/1	PodInitializing	0	17s	
init-demo	1/1	Running	0	18s.	Then it will go into Pod-initialization phase

Actual O' started.

> kubectl delete -f .

> kubectl apply -f .

> kubectl describe po init-demo

(<sup>Eg:-</sup> You shud never start app server without db server being up)  
we use init O'.

## ● Ephemeral Containers :-

These O's basically don't have any functionality.

If your pod is continuously failing → and if you want to debug (or) if you want to figure out what is the reason - why your O' is failing?

But, the point right now is → your O' is in exited state → You want to figure out the reasons → that is where Ephemeral O' come into play.



Google: Ephemeral Containers new feature of k8s (added in 1.25)

## Troubleshooting k8s

We will debug pods  
using  
Ephemeral CM

So,

This is for troubleshooting.

\* We know about Pods.

→ Pods we have 3 types of O's — 1) Normal O's

2) Init O<sub>rs</sub>

③ Ephemeral Ⓛ?

For every  $\sigma$  I can set the limits.

I can check whether my app is running or not by trying to write Probes.

## Node UseCases :-

\* Till now our focus was on Pods &  $\mathcal{O}^*$ s inside the Pod.

\* Let us assume that , you have Pods which run on your nodes .

You want to run a pod on a specific node. (till now we didn't worry about it).

In your org<sup>n</sup> there are some app<sup>s</sup> which require GPUs not CPUs.

Probably there might be some gaming server that is running inside the O which requires 11<sup>th</sup> processors - GPUs.

so in your k8s cluster → there are 2 types of nodes / one which have CPUs  
when you are running this particular pod, that needs to go and run on any node which has GPU -(not CPU) . → How do I do that ?

\* Remember Kube-scheduler

Till now we were trying to focus on Kubelet's functionality. How can I tell Kubelet that your  $\text{O}^*$  is not working. How can I tell Kubelet to restart the pod? This was what our focus was majorly on.

Now I'm trying to speak about Schedulers.

- 1) How to schedule a Pod on a Particular Node?
  - 2) How to stop assigning more Pods to a node?
  - 3) How to move all the pods running on a node to other node?

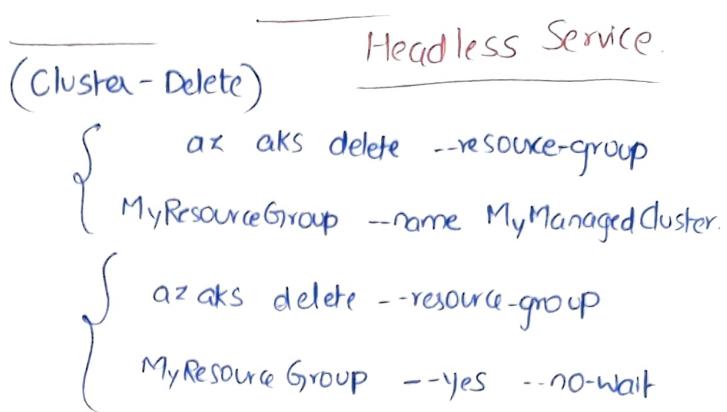
# Lec-57 - Apr-30 Part-1 | Deployment, DaemonSet, Annotations, Scheduling Pods,



## Deployment :-

→ Deployment creates a replicaset.

↓  
Creates Pod.



The basic thing which  $D^t$  can do is → **rollout** + **undorollout**

so, what do I mean by that? <----->

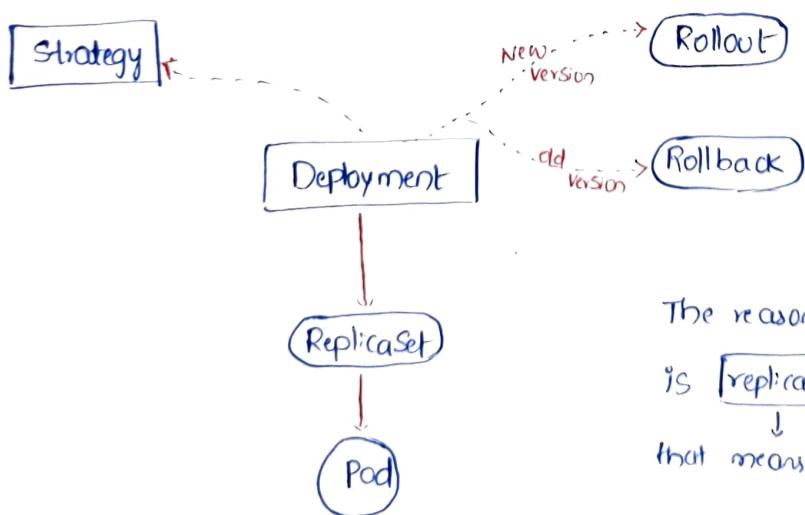
So, whenever a new version of your app comes → I can deploy.

when I deploy → there are strategies over here which can make your app  $D^t$  to be zero-downtime. → we can have a zero-downtime  $D^t$ . (zd- $D^t$ )

Now,  
if the new version of your app is not working → we can go back to the previous version.

→  $D^t$  is a k8s object which can help in rolling out and rolling back updates.

→  $D^t$  controls replicaset and replicaset controls Pods.



- Rollout is giving a new release

- Undo Rollout ≈ Rollback  
↓  
is giving a older version.

The reason why rs is diff. from rep controller is **replicaset can be persistent**

that means that →  $D^t$  tries to maintain your replicaset in etcd cluster

**Now, whenever you update, it stores the newer version of replicaset**

which is not possible using replication controllers.

\* So, there are two things which are different when it comes to replicaset and replication Controller :-

- 1) Replicaset can do Set-based Equality when it comes to pods.
- 2) Replicaset is used in  $\text{D}^t$  which is further used for persisting the previous versions.

So,  $\text{D}^t$  uses replicaset  $\rightarrow$  & in rs  $\rightarrow$  we try to write pods.

- One more imp thing in  $\text{D}^t$  is strategy.

So in  $\text{D}^t$ , this strategy tries to  $\rightarrow$  how do you want to do  $\text{D}^t$ ?

Do you want  $\text{D}^t$  to have downtime (or) do you want to have zdt?

↓

That is defined by this strategy.

$\rightarrow$  Apart from that, in  $\text{D}^t$  we are doing - basically calling replicaset and

Let's try this : Lets create a manifest with some application deployment :-

> mkdir deployments  $\rightarrow$  nginx-svc.yaml.

nginx-svc.yaml :-

---

apiVersion: v1

kind: Service

metadata

name: nginx-svc

spec:

selector:

app: nginx

ports:

- name: nginx-svc  
port: 80

targetPort: 80

Protocol: TCP

type: LoadBalancer

\* Now lets try to create → nginx-deploy.yaml  
Since if you try to go from one version of nginx to other version of nginx, you'll not be seeing much of a difference → it'll try to create a replicaset with image which is nginx → then it'll change the image to httpd → so that you can see the homepage difference.

→ API-reference → Deployment v1 apps

↓  
apiVersion → APPS/v1

Kind → Deployment.

name → nginx-deploy. (you can also add labels).

spec  
↓

Deployment Spec. → will contain <sup>a</sup> template → that template is ReplicaSet's template.

→ minReadySeconds = 1

replicas = 3

selector :

~~nginx~~.  
Type: LabelSelector → matchLabels  
app: nginx.

strategy

↓  
Type: DeploymentStrategy → RollingUpdate.

type: RollingUpdate.

rollingUpdate:->  
maxSurge  
maxUnavailable.

is default  
means it will delete the existing pod  
and then recreate the new one

whenever you put → your type: recreate  
→ you are having a down time.

How many replicas did we set? = 3.

For some time, let us assume that there are 3 pods which are of old version.

Now I'm trying to create a new version.

~~new version~~, whenever you are trying to tell strategy → You are trying to tell →

How do you want to do this deployment.

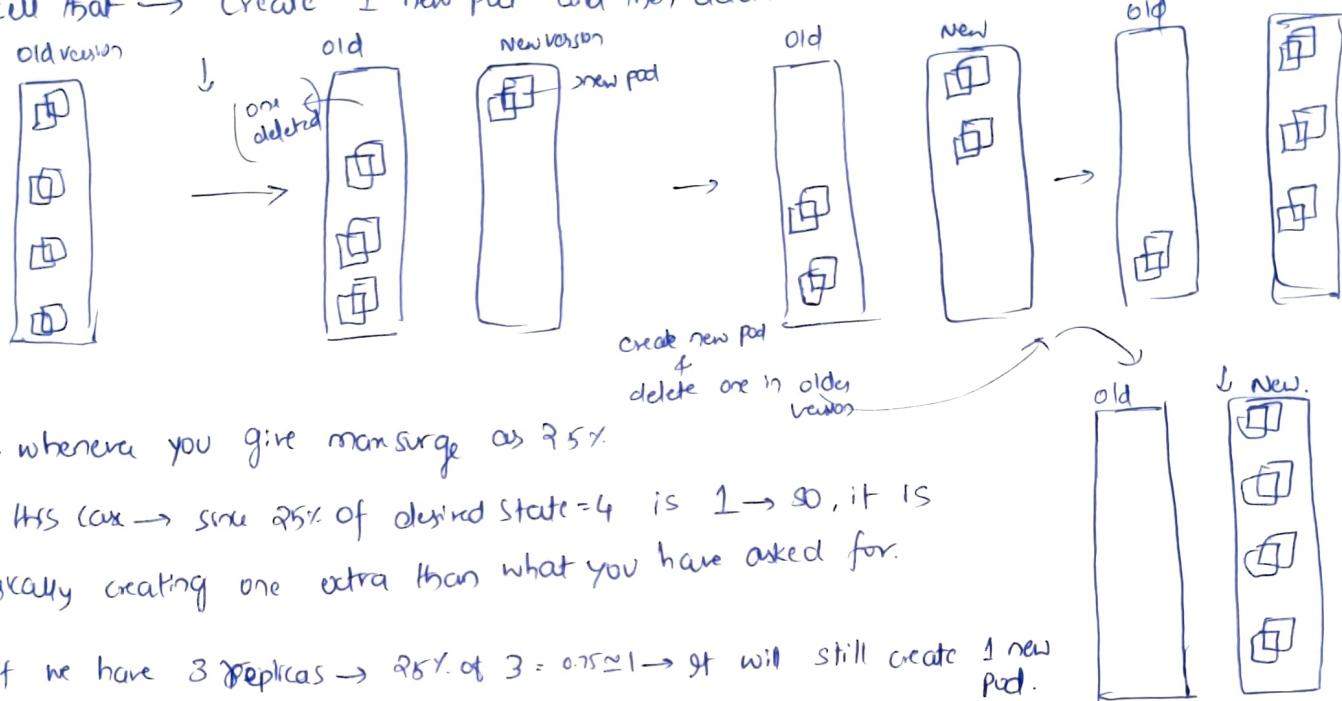
- So, basically we have 4 old pods (replicas=4)

NOW maxSurge  
↓ means

How many no. of pods (or % of pods which I can create above desired state).  
 $\hookrightarrow = 4$

Bcz I want new pods → to have new version.

So, let us assume that I have put maxSurge = 25%. → that means, I'm trying to tell that → Create 1 new pod. and then delete one here.



So, whenever you give maxSurge as 25%.

In this case → since 25% of desired state = 4 is  $1 \rightarrow$  so, it is basically creating one extra than what you have asked for.

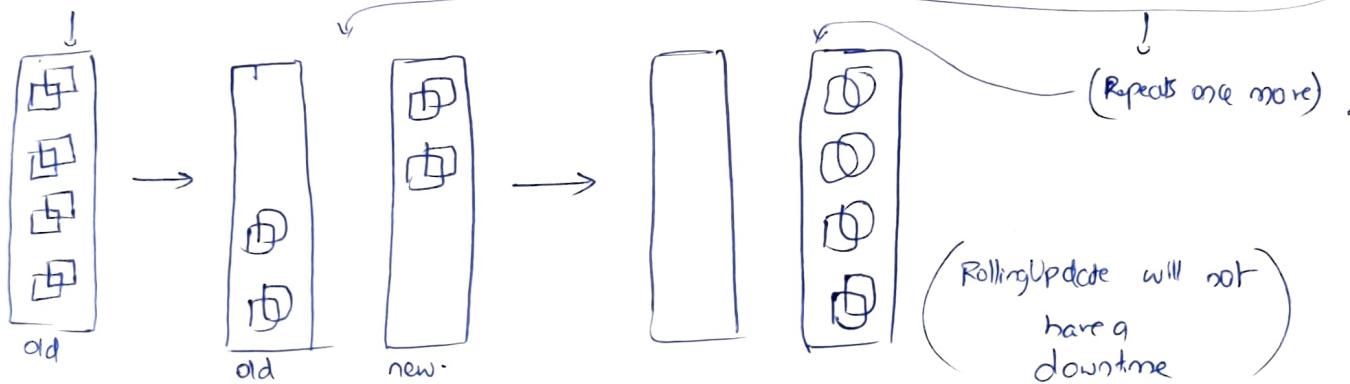
So, if we have 3 replicas → 25% of 3 =  $0.75 \approx 1 \rightarrow$  it will still create 1 new pod.

→ Let us assume → I have given maxSurge as 50%.

what it would do is

In this case → rather than creating 1 new Pod → it will create 2 at a time. Once those 2 are in running state, it will delete 2 in old.

bcz  $\rightarrow$  50% of desired state = 2.



If it is recreate → 1<sup>st</sup> it will delete the old Pod and then it will create a new one

→ that might have downtime. (bcz recreate is not one by one → it is everything at one shot)

(6)

## maxUnavailable :-

Whenever you are doing this  $\textcircled{D}^t \rightarrow$  how much % of pods can be unavailable.

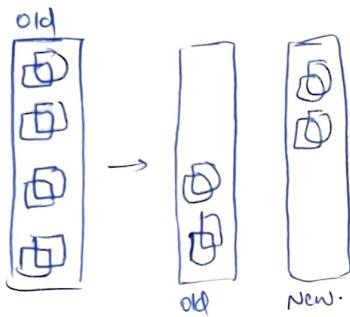
So, let us assume that

maxSurge  $\rightarrow$  I gave 50%.

(i.e. 2 pods at a time)

maxUnavailable  $\rightarrow$  I gave ~~25%~~ 25% = 1 pod. (Only 1 pod can be available)

so what it will do is



It will create 2 new pods & it will not delete 2 at one shot.

It will delete only one  $\rightarrow$  bcs only 1 pod can be unavailable and then it will delete 2<sup>nd</sup> pod.

That is how this  $\textcircled{D}^t$  strategy works.

maxSurge  $\rightarrow$  default no. is  $\rightarrow$  25%.

Here, you can write %. (or) numbers directly.

maxUnavailable  $\rightarrow$  25%.

$\rightarrow$  Next template (Pod template).

metadata :

: name: nginx

: labels:

: app: nginx

: ver: "1.23"

spec :

containers :

- name: nginx

image: "

ports :

- containerPort: 80

But I'm not using ver in SVC, as I might have 2 version.

so In my SVC  $\rightarrow$  I have only app:nginx - so I might get new version also. (so version can change).

## • nginx-deploy.yaml :-

---

```
apiVersion: apps/v1
```

```
kind: Deployment
```

```
metadata:
```

```
  name: nginx-deploy
```

```
spec:
```

```
  minReadySeconds: 1
```

```
  replicas: 3
```

```
  selector:
```

```
    matchLabels:
```

```
      app: nginx
```

```
  strategy:
```

```
    type: RollingUpdate
```

```
    rollingUpdate:
```

```
      maxSurge: 25%
```

```
      maxUnavailable: 25%
```

```
  template:
```

```
    metadata:
```

```
      name: nginx
```

```
      labels:
```

```
        app: nginx
```

```
        ver: "1.23"
```

```
  spec:
```

```
    containers:
```

```
      - name: nginx
```

```
      - image: nginx
```

```
    ports:
```

```
      - containerPort: 80
```

\* Google: kubectl cheat sheet (68)

kubectl deploy commands. (bluematador).

Exercise :-

Try writing deployment spec for a Jenkins app with 8 replicas and 50% surge + 50% unavailable.

Sp: jenkins/jenkins=lets-jdk11

\$ kubectl get nodes

2 nodes ✓

> cd deployments

> kubectl apply -f .

O/p: deployment.apps/nginx-deploy created

service/nginx-svc created.

> kubectl get svc

External-IP

O/p: Kubernetes  
nginx-svc LoadBalancer 10.0.114.92 20.232.81.56

Go chrome

http://20.232.81.56 → welcome to nginx!

\* Lets get deployment information:

\$ kubectl get deploy → short name.

O/p: nginx-deploy 3/3 3 3 . .

\$ kubectl describe deploy nginx-deploy

O/p: Name : nginx-deploy

Namespace : default

Annotations : deployment.kubernetes.io/revision: 1

Selector : app=nginx

Replicas : 3 desired | 3 updated | 3 total | 3 available | 0 unavailable

strategyType : RollingUpdate

RollingUpdateStrategy : 25% maxUn 125% max Surge.

\$ kubectl get rs Desired Current Ready

O/p: nginx-deploy-79...567 3 3 3

I have created a Deployment. It created an object.

\$ kubectl get po

↓  
Controls Pods

O/p: - - V<sub>1</sub> - -  
- - V<sub>1</sub> - -  
- - V<sub>1</sub>, - -

Lets explore rollout Command :-

\$ kubectl rollout --help

O/p: history → View rollout history  
pause  
restart → restart a resource.  
resume  
status → show status of rollout  
undo → Undo a previous rollout

\$ kubectl rollout history --help

\$ kubectl rollout history deployment/nginx-deploy

O/p Revision change-cause  
1 <none>

It says that you have only one version right now.

It has applied only once.

\$ kubectl rollout status deployment/nginx-deploy

O/p: deployment "nginx-deploy" successfully rolled out.

\* Lets update the specs to change image from nginx to httpd :-

Go to nginx-deployment.yaml → change image to httpd.

and minReadySeconds = 5 (for consideration)

That means - this Pod will come into running state after 5s.

> kubectl apply -f .

\$ kubectl get deployments.apps -w.

We can see that it has increased by 1 and

then it deletes the older version.

\$ kubectl rollout undo --help

You can just give current one.

But if you want to go to the Specific revision

Currently we have 2 revisions.

For some reason.

you have 100

if You want to go to 54

--to-revision=54

At the same time do

\$ kubectl rollout undo deployments/nginx-deploy

O/p: rolled back.

> kubectl get deploy -w

O/p: nginx-deploy

Name  
Ready  
3/3

Up-to-date  
3  
Available.  
3

" " 3/3 3 3

" " 3/3 3 3

" " 3/3 0 3

" " 3/3 0 3

" " 3/3 1 3

" " 4/3 1 3

4/3 1

" " 3/3 1

3/3 2

while this is happening ←

↓  
try to access

http://20.232.81.56

↓  
O/p: Welcome to nginx.

so, my app was — not down.

I was using my app continuously and then it moved to the older version.

\$ kubectl rollout history deployments/nginx-deploy

O/p: Revision      Change-cause

2 <none>

Now it is not showing 1.

3 <none>

Bcos 1 & 3 are same.

\* Now, let's try to create one more  $\textcircled{D}^+$  → image:- phpmyadmin.

\$ kubectl get deployments.apps -w

↓

O/p: = =

5/3

5/3

20.232.81.56 → Now it is Working fast as PHP does ~~not run~~ ~~run~~

maxSurge: 50%  
maxUnavailable: 50%

o/p: Name	Ready	Up-to-Date	Available	Age
nginx - deploy	4/3	2	3	8m 16s
nginx - deploy	4/3	2	4	8m 19s
" "	3/3	2	3	8m 19s
" "	3/3	3	3	8m 19s
" "	4/3	3	3	8m 20s
" "	4/3	3	4	8m 25s

In the meantime

In chrome

httpd  
http://20.232.81.56 → It works!

So, am I able to access the app? while  $\textcircled{D}^t$  is happening?

YES.

\$ kubectl get rs

o/p: nginx-deploy-787...959  
nginx-deploy-79d...567

} 2 replica sets

one refers to previous version  
when we ( $1^{\text{st}} \textcircled{D}^t$ )  
have done ↑

$2^{\text{nd}}$  one → since I gave new version  
it created new rs.

\$ kubectl get po

o/p: newer versions of pods.

\$ kubectl rollout history deployments/nginx-deploy

o/p: Revision Change-Cause

1	<none>	→ $1^{\text{st}}$ one was with nginx
2	<none>	2 $^{\text{nd}}$ " " " httpd.

\$ kubectl rollout status deployments/nginx-deploy

if for some reason → let us assume that → you have given a new version → and it is not working acc. to your customer's expectations → then what we would want to do is → we would want to go back to the previous version.

→ Now To roll back to Previous Versions. ←

\$ kubectl rollout --help.

o/p: Undo.

It is failing. So what should we do? → Go back to previous version.  
working.

(70)

> kubectl rollout history deployment/nginx-deploy

Revision      change-cause

~~000b~~ 2      <none>

0/p: 3      "

4 → only name (php) - sir's mistake.

5 → image also

> kubectl rollout undo deployment nginx-deploy --to-revision=2

0/p: rolled-back.

So I'm able to move from one version to other.

> kubectl rollout history deployment nginx-deploy

0/p: 3

4

5

6

what is change-cause?

It is showing as none which is not good.

while applying

what can be done to have a valid change-cause?

So, change-cause could be as simple as → I have updated to this new revision bcoz

(or)

just write some info.

↓ I want to write some info

I got a new build on monday

↑ on just write date here.

For that I need to do certain things. → which could be → making some arguments to be passed while using kubectl apply command.

## ● Annotations :-

We spoke about Labels. Labels are → which are used by k8s for selections. There is one more Key-value Pair in k8s → but this is not meant for k8s → this is meant for tools which are using k8s → to add some logic.

→ For Eg! - Yesterday we created a LBalancer. Now, I want to create a LB in Azure but I don't want public loadbal. → i.e. I don't want it to create a public IP add. I want it to create a private-IP which works only within the network of Azure. So, how do I tell that information?

- so k8s manifest file doesn't have any such info.

so, Generally, those tools use something called as Annotations.

→ Annotations are key-value pairs which are used by External tools.

For Eg: Anything specific to Azure → if you want to write in your manifest → you can't use any other sections → the only section which is left to you is → Annotations.

Google: annotations in k8s.

● Annotations are also <sup>like</sup> metadata. But they will be used by tools.

So, the way you write annotations is not very much different from the way you write Labels. But, the usecase is different.

Google: k8s ⚡ Ann's. | Change-cause Annotation.

whenever you want to basically deploy → what you can try to write is → You can try to write an Ann' which is called as change-cause . ✓

→ Ann's have command line also. (just like kubectl label) ., but it is not a good idea to use commands. Try to write them in Deployment itself .

# Google: k8s Service annotations azure (Microsoft)

↓  
use an Internal Load balancer with AKS

so, creating an Internal lb in Azure.

Service → type: LoadBalancer

metadata

Annotations:

service.beta.kubernetes.io/azure-load-balancer-internal: "true"

we are telling Azure to do this  
so CCM will take control  
and do it.  
K8s has not control over  
it

So, the moment you add this ann? → what Azure will do is → OK, he is asking for a lb but it shud not be public.

- So, for every cloud provider → You wud have some specific annotations,

Especially whenever we are dealing with ingress and some cloudspecific things.

In those cases, we wud be writing ann?

→ Go to nginx-deploy.yaml

↓  
apiVersion: apps/v1

kind: Deployment

metadata:

name: nginx-deploy

Annotations:

kubernetes.io/change-cause: "update to phpmyadmin"

(image: Phpmyadmin)

> kubectl apply -f .

> kubectl rollout history deployment nginx-deploy

Op:	Revision	CHANGE-CAUSE
	3	none
	4	"
	5	"
	6	"
	7	update to phpmyadmin.

so, whenever you change a deployment,  
change that annotation and write  
reason why you are doing that.

This is the better way.

Google: k8s Azure Annotations

Google: Aws k8s Annotation Specs.

(Github)

Aws k8s Service Annotations (Service related aspects)

The Usecase of Annotations is → to add any external functionality by third-party tools. This is not something K8s will use internally

Google: k8s Service annotations azure (Microsoft)

↳ Use an Internal Load balancer with AKS

↳ So, creating an Internal lb in Azure.

Source → type: LoadBalancer

↳ metadata

Annotations:

service.beta.kubernetes.io/azure-load-balancer-internal: "true"

we are telling Azure to do this  
so ccm will take control  
& do it.  
K8s has not control over  
it

So, the moment you add this ann? → what azure will do is → OK, he is asking for a lb but it shud not be public.

• So, for every cloud provider → You wud have some specific annotations,

Especially whenever we are dealing with ingress and some cloudspecific things.

In those cases, we wud be writing ann's.

→ Go to nginx-deploy.yaml



apiVersion: apps/v1

kind: Deployment

metadata:

name: nginx-deploy

Annotations:

kubernetes.io/change-cause: "update to phpmyadmin"

(image: Phpmyadmin)

> kubectl apply -f .

> kubectl rollout history deployment nginx-deploy

Op: Revision CHANGE-Cause

3 none

4 "

5 "

6 "

7 update to phpmyadmin.

so, whenever you change a deployment,  
change that annotation and write  
reason why you are doing that.

This is the better way.

↳ Google k8s Azure Annotations

Google: Aws k8s Annotation Specs.

(Github).

Aws k8s Service Annotations (Service related aspects)

The Use case of Annotations is → to add any external functionality by third-party tools.

This is not something K8s will use internally

## ● DaemonSet :-

\* DaemonSet is a controller which creates pod on every/selected nodes in K8s cluster.

For Eg:- if you want to run one pod for every node in your cluster.

Use Cases :-

- log collectors
- agents etc.

This is not an application. This is additional functionality.

For Eg:- In K8s nodes, you have certain data → You want to export that data all the time from every node in your K8s cluster → so, how do you do that?

You basically create a DaemonSet.

DaemonSet ensures that whatever Pod you have configured in Daemonset runs on every node in K8s cluster.

Basic use case is → whenever you want to run agent-kind of ~~things~~ (Agents are the softwares that get installed on every node and then they communicate with server)

→ Let us assume that I want to install some agent → log agent → which collects logs and exports to Server.

Any kind of agenting kind of functionality → we will be using a daemonset.

Google: K8s daemonset.

→ This is again a controller which controls pods.

when you delete the node → pod will automatically be deleted. It will not try to reassign-bcoz it is speaking about one pod per node.

Some typical ~~usecases~~ of a Daemonset are :-

- running a cluster storage daemon on every node.
- running a logs collection daemon on every node.
- running a node monitoring daemon on every node.

> mkdir daemonset

\$ kubectl api-resources | grep daemonset



fluentd-ds.yaml :-



ds → is short form.

fluentd is one app which is log collector. It can read logs and it can send to whatever server you have configured.

so we are configuring fluentd.

You can configure different kinds of logs.

Google: api reference → Daemonset v1 apps

apiVersion: apps/v1.

kind: DaemonSet

name: fluentd-ds

spec annotations: kubernetes.io/change-cause: "added fluentd debion".  
↓

DaemonSet spec → minReadySeconds: 5

Here also, you can do deployments, rollout, underrollout - - -

revisionHistoryLimit → how many old history do you want to retain. to allow rollback

default = 10. (Previous 10 things)

selector → LabelSelector:

matchLabels

app: fluentd.

template → PodTemplateSpec.

metadata:

name: fluentd

Labels:

app: fluentd.

spec:

containers:

→ name: fluentd

image: fluentd:v1.16-debian-1.

Lets try to run fluentd container

Playground.

docker hub: fluentd.

\$ docker container run -d -P fluentd

\$ docker container ls



## fluentd-ds.yaml :-

```
apiVersion: apps/v1
Kind: DaemonSet
metadata:
  name: fluentd-ds
  annotations:
    kubernetes.io/change-cause: "update to latest"      ↑      1st case
spec:
  minReadySeconds: 5
  selector:
    matchLabels:
      app: fluentd
  template:
    metadata:
      name: fluentd
      labels:
        app: fluentd
    spec:
      containers:
        - name: fluentd
          image: fluentd:latest.                         fluentd:v1.16-debian-1
```

> kubectl apply -f .

O/p: daemonset.apps/fluentd-ds created

> kubectl get ds

	Name	Desired	Current	Ready	Available	Node Selector
	fluentd-ds	2	2	0	0	<none>

I have 2 nodes and it is running on both the nodes. I have not specified any count. The only thing which I specified is - I want a Pod called fluentd.

> kubectl get nodes

I have 2 nodes. It is running on both.

→ If I basically increase my count → to 3 nodes → then it will automatically create 3 pods. That is how daemon sets work.

\* > kubectl rollout history ds fluentd-ds

<u>Revision</u>	<u>Change-Cause</u>
1	added fluentd debian

> kubectl apply -f .

> kubectl rollout history ds fluentd-ds

<u>op</u>	<u>Revision</u>	<u>Change-Cause</u>
	1	added fluentd debian
	2	update to latest

So  
Let it be daemonsets, let it be deployments, let it be stateful sets  
Every time add that annotations.

→ For 3 controllers → we can see rollout status — 1) Deployments  
2) Daemonsets &  
3) Stateful sets.

\$ kubectl get po

o/p: fluentd-ds-kk. 1/1 Running  
fluentd-ds-hpc. 1/1 Running

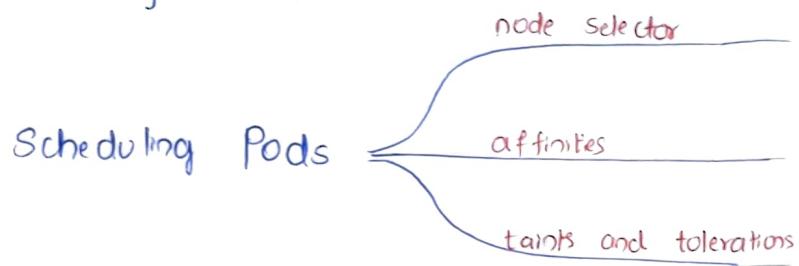
\* kubectl delete -f .

\* I want to schedule pod on some nodes.

So what are the options which k8s will try to give me.

## ● Scheduling Pods :-

There are 3 major ways :-



→ One is called as node selectors.

In node Selectors, what you would try to tell is,

in your podspec you will directly give node name (In which node you want to run the Pod)

→ Affinities and taints + tolerations are a bit tricky. Node selector is easy method.

### ## Node Selectors :-

> mkdir scheduling

> cd scheduling

> mkdir nodeSelector.



jenkins.yaml

> kubectl get nodes

O/p: aks-nodepool1-000  
aks-..-001  
I have 2 nodes.

I want to run it on node-000 (whatever Pod which I'm trying to write).

→ Google: node selector k8s (Assigning Pods to nodes)

Docker hub: jenkins → jenkins/jenkins:jdk11

## jenkins.yaml

apiVersion: v1

Kind: Pod

metadata:

name: nodeSelector  
labels:

spec:

labels:

app: nginx

purpose: nodeSelector

containers:

- name: jenkins

image: jenkins/jenkins:jdk11

ports:

- containerPort: 8080

## expose.yaml

---

apiVersion: v1

kind: Service

metadata:

name: expose-svc

Spec:

selector:

purpose: nodeSelector

type: LoadBalancer

ports:

- name: nginx-svc

port: 80 → (if I give 8080 in port)

targetPort: 8080 ↓ then I have to give loadbalancer's IP: 8080

Protocol: TCP while accessing Jenkins.

> kubectl apply -f \nodeSelector → folder's name (whatever is present in the folder)

\$ kubectl get po -o wide will be created)

Name	Status	Node	Nominated Node	Readiness Grade
nodeSelector	Running	aks-node-001	<none>	<none>

- so where did the pod get created? → on node-001.
- \* To be very specific → I didn't write any logic (ie. I didn't give which node I want my pod on). It might have been created anywhere.
  - \* So if I create one more pod → it might get created on node-000 (or) node-001. Right now, I'm not controlling → on which M/c(node) my (pod) app<sup>n</sup> is running

If I want to access my app!—

">\$ kubectl get svc

o/p: exposc-svc LoadBalancer -- 20.75.165.138 80:30890/TCP

Kubernetes clusterIP --

Chrome: → 20.75.165.138 → open.

\* Now, I'm interested in which node my app runs on:—

Go to API ref → Pod v1 Core.

↓  
Spec → nodeSelector, nodeName.

↓  
object.

\* Now lets Select node by its name:

↓  
But it is not a good idea.

> kubectl get nodes

o/p: Name  
aks-nodepool1...000 → I want my app<sup>n</sup> on this node.

aks-nodepool1...001

First → Edit jenkins.yaml

↓  
nodeName: aks-nodepool...000. (Is this a wonderful way?)

Spec:

nodeName: " "

containers:

- name: jenkins

↓  
No,

will it work? → Yes.

> Kubectly apply -f .

\$ kubectl get po -o wide

A.S.

{ Since the pod was already created  
it is not a controller object  
You can't do changes .

Name	Nominated Node	Node
nodeselector	1/1	aks-000 ✓

Now, our pod is running on whatever node we specified in yaml file.

Using nodeSelector (or) any others → to assign pods to nodes is not good idea.

\* Rather than using this,

let us assume that, I have categories of nodes.

So, how can we select specific pods in k8s → we will go with Labels.

So, if you apply the same logic → so, while I'm basically creating nodes, for the nodes I will either use existing labels (or) I'll create labels and rely on those labels → bcoz →

let us assume that - I have a label called as hardware type : GPU →

so what I'm trying to tell in node selectors is → I would want to run my pod on any of the nodes which have this label.

This is a better way rather than giving a node name (or) node's ip-address bcoz we will never know which node is available or which node is not available. So relying on labels might be a good idea.

\* Lets look at nodeSelectors :-

→ Node Selector is a selector which must be true for the pod to fit on a node.

Selector which must match a node's label for the pod to be scheduled on that node.

→ We have two nodes. Lets attach the following labels:-

- Purpose: poc — node 0
- Purpose: testing — node 1.

For poc purposes, if I'm developing a pod then I want to use node 0.

For testing whatever I have developed, I want to use node 1.

\$ kubectl label nodes aks-nodepool1-328c-0000 purpose=POC

Output: Error.

There is no other way. I have to add another nodepool.

Our nodepool's name is nodepool1.

\* Go to my AKS cluster → nodepool1 →

nodepools → scale nodepool



change node count from 2 to 1.

I'm just making count back to 1.

\* I'll try to add one more nodepool: nodepool2

\$ az aks nodepool add \

--resource-group myResourceGroup \

--cluster-name myAKSCluster \

--name nodepool2

--node-count 1 \

--labels dept=HR purpose=fusting

--no-wait.

ideally you can apply labels directly.

But, since this is a managed cluster it is not allowing you to do that.

az aks nodepool update \

--resource-group myResourceGroup \

--cluster-name myAKSCluster \

--name labelnp \

--labels dept=ACCT purpose=POC

--no-wait

↓  
see for labels in Azure.

> kubectl get nodes

Output: aks-nodepool1-000 Ready

aks-nodepool1-001 Ready/scheduling Disabled.

Eventually  
we can see other nodepool coming up.

↓  
is disabled. (later discussion)

Now, go to Jenkins.yaml:-

replace nodename with nodeSelector.

Spec

nodeSelector:

purpose=POC

label which I want to

match is

Any node which has this label → ... will run this Pod.

\$ kubectl get nodes

O/p: aks-nodepool1-000 Ready ---

Go to Azure → my AKS cluster → we can see nodepool1 - succeeded  
nodepool2 - creating.

what we want is not this ↙

> kubectl apply -f ./nodeSelector

> kubectl delete -f ./nodeSelector

> kubectl apply -f ./expose.yaml

> kubectl apply -f ./jenkins.yaml

> kubectl get po -o wide

It is going to nodepool1 bcoz of label matching.

Let us see if second nodepool is up.

\$ kubectl get nodes

O/p: aks-nodepool1-000 ✓

aks-nodepool2-000 ✓

> kubectl delete ./jenkins.yaml

Now edit jenkins.yaml → purpose = testing. → we have given this label in nodepool2.

> kubectl apply -f ./jenkins.yaml

> kubectl get po -o wide

O/p: nodeSelector --- aks-nodepool2 . It went to nodepool2.

Kubernetes.io/os: linux

Kubernetes.io/os: windows

Now, you would write a node selector where you'll tell that - I'm interested in running this pod with . Then your pod will end up on Windows Server if they exist.

> kubectl delete -f \jenkins.yaml

This time let's give a purpose value which is not existing.

Jenkins.yaml → purpose: development. ( purpose: poc - nodepool1  
purpose: testing - nodepool2 )

> kubectl apply -f \jenkins.yaml.

↓  
we have only 2 labels.

> kubectl get po -o wide

	Ready	Status	Node	Nominated Node
0/p: node selector	0/1	Pending	<none>	<none>

Your k8s Pod is in Pending state. <sup>Bcoz</sup> Your Scheduler is unable to find any node which has a label called as purpose: development.

So Kube-Scheduler is unable to decide which node your pod should be send to. When label doesn't match → your pod will always be in pending state.

→ when we have tried to create a Pod with node selector matching purpose:poc it was created on nodepool1 and

when we created a pod with purpose:testing it created in nodepool2 and

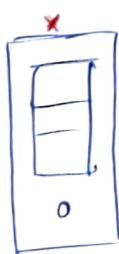
when created a pod with purpose:development it was created in Pending state (not created)

\* when I want to run a pod → I want to run a pod in a node which has these conditions → that is node affinity.

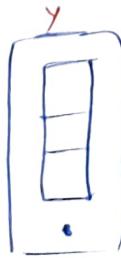
I want to run a pod in any node where there are pods with these settings

→ that is Pod affinity.

Let's assume I have 2 servers. Now I have to schedule a Pod.



In this there is already a pod called as 'A'  
(A + B are conditions)



In this there is already → Pod 'B'.  
(A + B are conditions)

I can try to tell that now create a Pod based on condition X.

X & Y are node related conditions,

A + B are pod related conditions.

then it will run in 1st one.

- So, whenever you are relying on node-based conditions → that is 77  
Called as node affinity
  - Now, if I try to tell that → I want to run this pod where there is no 'B' running  
then it will go & run on 1<sup>st</sup> one.  
I want to run this pod where 'B' is running → then it will go & run on 2<sup>nd</sup> one.  
This is called as Pod affinity.
- Here we are trying to help Kube-Scheduler to find the right node to run your workload.

In an ideal case → You should not worry about it.

The reason why we basically do these kinds of things is w.r.t diff hardware  
So, I might have nodepool1 where there is 4GB Ram node.

I might .. nodepool2 .. .. 16GB .. ..

This pod for some reason has to run on 16GB node. So then how do I schedule it?  
This pod has to run on node where there are 8 CPUs → .. ..

\* I don't want to run 2 mysql pods on the same node — for some reason,  
↓  
Then I'll try to go with Pod affinities.

→ There is something called as affinity and anti-affinity.

### Affinity/Anti-Affinity Based :-

\* This scheduling happens O.b.O affinity.

Affinity means when it matches schedule it here.

Anti-Affinity means when it matches → don't run it here.

Google: node affinity in K8S → Assign pods to nodes using Node Affinity.

Here you will have 2 interesting sections

In Pod → nodeAffinity:  
→ requiredDuringSchedulingIgnoredDuringExecution

• This basically means that

when you are using required section

whenever your pod is getting scheduled - this shud match →

when it matches only — your pod will be selected

But when you are using preferred section :-

nodeAffinity:

PreferredDuringSchedulingIgnoredDuringExecution

It will try to see if it matches → Even if it doesn't match also → It will schedule.

★ Preferred is not a mandatory rule → It is 'good to have' kind of rule.

Required means → it shud mandatorily follow.

In earlier case → when labels were not matching → Your Pod was in Pending state.

In this case, what wud happen is

So, if you use it in preferred section → Even if it is not matching → it wud apply somewhere → It will create a Pod.

API-ref

→ Podspec



affinity (it is basically matching constraints)

Type(Affinity)

mandatory rule

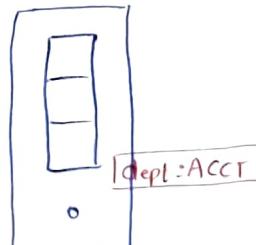
↳ nodeAffinity → Speaks about labels of the node, good to have rule.

Pod Affinity → Speaks about labels of other Pods which are running on that node.

Pod AntiAffinity

↳ means I want to run this pod on a node which doesn't have Pod with this label.

So, there are 3 kinds of rules over here.



Now, whenever I'm trying to write an affinity rule.

↓  
I'm trying to write an affinity rule which is based on node affinity

↓  
I'll try to tell → dept:HR [House Icon]

(But I have written this in requiredDuring...)

So then what wud happen is

↓  
It will try to see that, there is a label called dept:HR in 1st node  
→ Let me try to run it over here.

\* Now, I have written this label not in 'required' but in 'preferred'

↓  
So, there is no node (not necessary to match) (Dept: HT)

that is matching it

↓

But this might be executed on any one of the node.

→ This is how node-affinities work.

\* Now, let's speak about Pod-Affinity :-

I'm trying to apply Pod affinity which is trying to tell that → [app: jenkins]

so, now, what I'm trying to tell is → I want to run this pod on a node which has any pod → which has a label → [app: jenkins] → Now it will go & schedule on 2nd Server(node).

\* Now, if I have written this in Anti-Affinity Section →

Since there is already a Pod with label → [app: jenkins] → It will not schedule over here (in 2nd node)  
→ It will schedule it over 1st node.

\* Taints and Tolerations :-

Google: Taints and Tolerations in K8S.

Node affinity is basically the property where you are trying to tell that → run this pod on a node where there are matching labels.

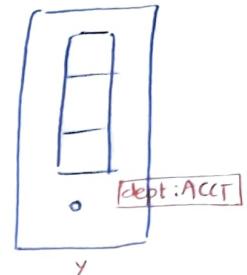
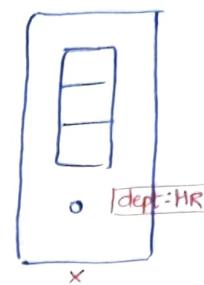
But here, it is exactly Opposite.

Taints, what they try to do is → They allow node to → not allow pod to run over here.

→ Taints are the opposite -- they allow a node to repel a set of pods.

• So, basically what you would do is → You would apply taint to <sup>the</sup> node.

• Taint is some condition.



Now, there are tolerations.

Tolerations is something which we apply to the pod.

\* Now k8s will apply pod to a node with matching tolerations.

→ In all the cases, what is it that we are trying to do →

I have a pod → I want it assigned to the node → Here, based on the Podspec — we are making a decision.

But, over here, the case is reverse → Here, node has a taint and based on the taint → pod will be scheduled.

Eg: tainting a node → kubectl taint nodes node1 key1=value1:NoSchedule.

↓  
places a taint on node → (node1).

The taint has Key=Key1 and Value=value1 and taint effect → NoSchedule. This means that no pod will be able to schedule onto node1 unless it has a matching toleration.

Spec:

Eg: tolerations:

- Key: "key1"  
operator: "Equal"  
value: "value1"  
effect: "NoSchedule"

tolerations:

- Key: "key1"  
operator: "Exists"  
effect: "NoSchedule"

Taints & tolerations → because

↓  
These things are reqd  
when you have nodes  
with specialized hardware

↓  
if key1 exists do not schedule.

It is written

that → when key1 = value1 → do not schedule.

The default value for operator is Equal

\* Here, the priority is given to the node and whatever taint is applied to the node → there should be a toleration against it.

Pod-Affinity, node selector → these are all priorities given to the pod and node has no decision to make.

But here → a decision-maker is actually a node

(So, in the node, we will be trying to apply a taint saying that → apply pods only when they have matching tolerations, otherwise do not schedule)

- k8s is a cluster.

So let us assume that you have created a AKS cluster. And now, your Org<sup>n</sup> has 5 or 6 Environments. Dev, QA1, QA2, UAT, staging, prod? Envs.

So, do I need to create 6 K8S clusters for that. The answer would be Yes (or) no.

But the point is, k8s tries to tell that → I have a virtual cluster in me. You can create a virtual cluster, not a physical cluster. Virtual cluster → which is called as namespace.

All these we are using namespace and our namespace is default.

## • Headless Service in k8s :-

Google → k8s headless Service → Service

Generally whenever we create a Service → Every Service gets an IP-address.

```
mkdir headless-svc
```

## nginx-svc.yaml :-

apiService: v1

Kind : Service

metadata :

name = nginx-sv

spec:

selector:

app: nginx

type : ClusterIP

## Cluster

orts:

name: ngi

port: 80

Target Port: 80  
Protocol: TCP

I'm speaking about type: clusterIP.

when we create clusterIP → we get IP address

ClustalP: none

↓ what this means is → I will not  
the Service -

Any service which does not have its own IP-address is called as a Headless Service.

## nginx-rs.yaml :-

```

apiVersion: apps/v1
kind: ReplicaSet
metadata:
  name: nginx-rs
  labels:
    app: nginx
    layer: web
    version: "1.23"
spec:
  minReadySeconds: 1
  replicas: 5
  selector:
    matchLabels:
      app: nginx
  template:
    metadata:
      name: nginx
      labels:
        app: nginx
        version: "1.23"

```

## Spec :

```

    containers:
      - name: nginx
        image: nginx:1.23
      ports:
        - containerPort: 80

```

> kubectl apply -f .

%p: pod/alpine created

replicaset.apps/nginx-rs created

service/nginx-svc created.

## alpine-pod.yaml :-

```

apiVersion: v1
kind: Pod
metadata:
  name: alpine
spec:
  containers:
    - name: alpine
      image: alpine
      command:
        - sleep
        - 1d

```

> kubectl get po

o/p: 1 alpine

5 replicas.

> kubectl get svc

	External IP	
kubernetes	<none>	443/TCP
nginx-svc	clusterIP	80/TCP

\$ kubectl exec -it alpine -- /bin/sh

/# nslookup nginx-svc

O/p: I got 5 IP-addresses of k8s replica sets.

When you query a Headless Service → it will not give its IP-addr. It will give IP-addresses of all the pods which (it is connected to)  
(<sup>(or)</sup> it exposes)

This is useful when we try to learn something called as Stateful Sets.

But if write cluster-IP → I'll not be getting IP-addresses of Pods. I'll just get Service's clusterIP's IP-address.

\* W.r.t storage, first of all, let us quickly reiterate the classical problem with Docker Containers.

The problem with DR<sup>o's</sup> is → Every DR<sup>r</sup> whenever it writes any content into the local system → it would be present in Read-write Layer and when you delete the DR<sup>r</sup> → you will lose all that data.

→ To resolve this, in the case of DR<sup>r</sup>, we have used Volumes.

The basic idea is, a volume which we have created (o) a which helps in retaining data → how can I use it in K8S.

what are the different options that are basically present?

- Stateful Applications store data locally. In Containers the data created locally will be lost once you delete it. So to solve this in docker, we have used volumes. Volumes have a lifecycle which has no relation to Container Lifecycle (ref a DR<sup>r</sup>, image layers, volumes)

→ For eg: If I have a Web app<sup>n</sup> and whenever I try to access the web app<sup>n</sup> & if it is storing all the data in database → then this web app<sup>n</sup> will become stateful and that database becomes stateful. ("stateful app<sup>n</sup> use local storage to store some data whereas stateless app<sup>n</sup>s do not do that")

\* In K8S, we are running docker containers, K8S is an orchestration solution.

Lets see what are options for storage Provisioning in K8S

Google: K8S storage.

→ The most widely used storage types :-

- Volumes
- Persistent Volumes
  - Storage classes
  - Persistent Volume Claims.

\* All the other storage options you see in k8s are basically an extension to the ones mentioned before.

## ● Volumes :-

A Pod can have any number of volumes. But this volume will be available as long as Pod is running. Once you delete the Pod → volume will also be lost.

But in the Pod, no matter how many times your ⚡ will restart → it will still have the same volume.

The scope of this volume is — the lifetime of Pod. Once you delete the Pod you will lose the volume,

still, this volume will have lifetime other than container. No matter how many times ⚡ restarts → volume is preserved as long as Pod is alive.

● Volumes can be mounted to containers and they have lifetime equivalent to Pods.

→ It is not a very useful scenario for us → bcoz, in k8s we generally don't create ⚡'s directly → we create pods → so, when you are deleting the pods, you are losing the volume.

→ Types of volumes :

\* storage on cloud

- ebs
- azure disk
- efs
- azure file
- gcs

\* Empty dir → means when you create a volume → it will give you an empty directory.

\* hostPath → means → your Pod runs on some node. On the node there is some folder (or file) which will be mounted.

\* Lets create a manifest with mysql pod with volume

\* Manifests → mkdir storage → mkdir volumes → mysql-vol.yaml

mysql-vol.yaml :-

---

apiVersion: v1

kind: Pod

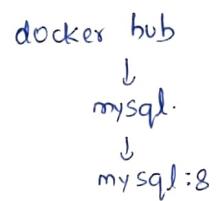
metadata:

name: mysql-vol

spec:  
labels:  
app: mysql  
layer: db

containers:

- name: mysql
- image: mysql:8
- ports:
- containerPort: 3306



volumeMounts:

- name: test-volume
- mountPath: /var/lib/mysql

volumes:

- name: test-volume

emptyDir:

sizeLimit: 100Mi

In Q, whenever I want to try to use a volume,  
I will use → volumeMount.

↓  
array.  
will have mountPath

↓  
name  
must match with  
name given in Volume

↓  
that is how the link  
gets created

volume ↴

This is not useful.

we have something that is useful → Persistent volume.

Persistent Volumes :-

→ These volumes will have a lifetime different than a Pod → i.e. they exist even after the Pod is dead.

→ Types of PVC (Google: K8s Persistent Volumes)

→ Storage classes in Azure (Google: AKS Storage classes)

→ Storage classes in AWS (Google: AWS Storage classes)

There are a kinds of Persistent volumes:-

① Static Provisioning → Your K8S administrator needs to create the storage.

② Dynamic Provisioning

↳ when you go with Dynamic volume,

K8S will automatically create the storage dynamically  
by taking the help of CCM.

In Azure → it will create Azure storage

In AWS → .. .. AWS ..

• So, we are speaking about persistent volumes which

not only have lifetime other than your pods, they also can create storage for you.

• For creating the storage, what they use is, they use something called as Storage classes.

In many of the managed K8S clusters → storage classes already exist.

So, you can take the help of them to create this dynamic storage (or) you can write something on their own.

• So, persistent volume is a disk.

Storage class is something which can help you in creating disk.

For Eg: If I want 1 GB of storage → so creating that 1 GB of storage is something which storage class can do.

• In your pod what we can do is → we make a Persistent Volume claim (PVC)

So, whenever we are trying to write the pod → we'll be trying to create a PVC.

• PVC :-

\* A PVC is a request for storage by a user. It is similar to a Pod.

Pods consume node resources and PVCs consume PV resources.

\* Pods can request specific levels of resources (CPU and Memory).

Claims can request specific size and access modes (Eg. They can be mounted ReadWriteOnce, ReadOnlyMany (or) ReadWriteMany).

- Persistent volume is a actual storage which can be dynamically created with the help of storage classes,

PVC is a link b/w your Pod and the volume that is actually created.

\* We have a kind → PersistentVolumeClaim.

Kind → PersistentVolume

In your pod, whenever you want to use the mount, you will be using this claim.

→ Types of PVC :-

\* Azure disk

\* azure file

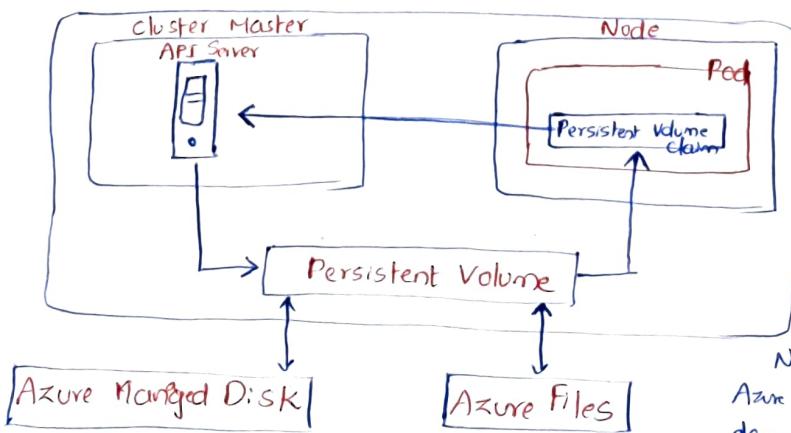
\* AWS EBS

\* gce Persistent Disk etc.

\* CSI → <sup>is a</sup> Container Storage Interface → which can be implemented by any plugin  
→ thus they can create a storage for you.

For Eg: If your company is using NetApp storage → there are companies who sell NetApp's CSI-plugin. Now you need to add it to the k8s cluster → Once you do that → Your manifest will be able to create NetApp's storage.

\* Google: Azure AKS persistent volume. AKS Cluster



whenever you are creating a node,

In the node, you have a pod

In the pod, you create PVC.

This PVC will basically get attached/linked to your Pod with the help of Persistent Volume.

Now whether you want to create Azure Managed disk or Azure file depends on storage class.

\* You can use different kinds of storages.

→ Volumes defined and created as part of the Pod lifecycle only exist until you delete the Pod.

Pod often expect their storage to remain if a Pod is rescheduled on a different host during a maintenance event, especially in stateful sets.

→ A Persistent Volume (PV) is a storage resource created and managed by the K8S API that can exist beyond the lifetime of an individual Pod.

→ You can use Azure Disk or Azure files to provide the P.V.

- To define different tiers of storage (such as premium & standard), you can create a StorageClass.

The StorageClass also defines the reclaimPolicy.

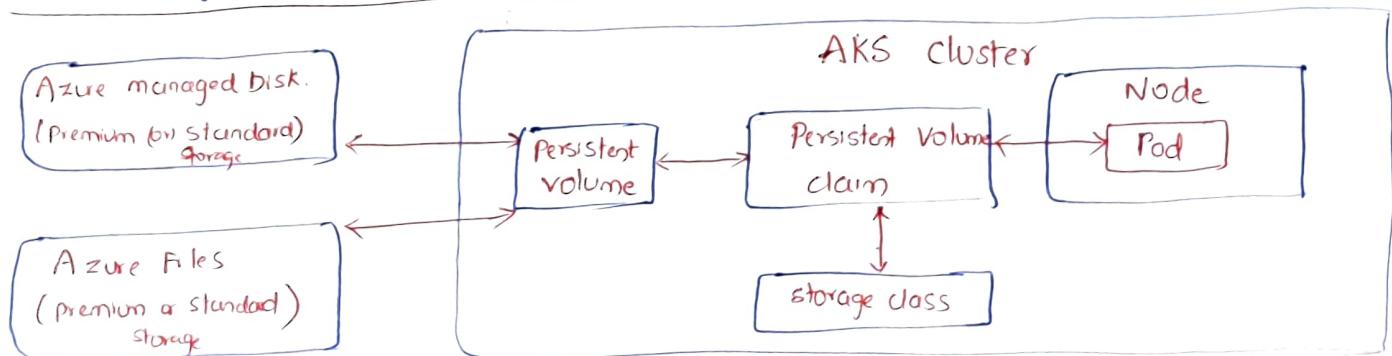
When you delete the P.V, the reclaimPolicy controls the behaviour of the underlying Azure storage resource. The underlying storage resource can either be deleted or kept for use with a future Pod.

- Storage Classes :-
  - use managed-csi → to get standard SSD (you get standard SSD disk)
  - managed-csi premium → creates a Azure premiumdisk.
  - Azurefile-csi-premium → creates a Azure file share.

\* If you want, you can create a storage class for additional needs using Kubectl.

→ A PVC requests storage of a particular storage class, access mode and size.

The K8S API Server can dynamically provision the underlying Azure storage resource if no existing resource can fulfill the claim based on the defined StorageClass.



- Once an available storage resource has been assigned to the Pod requesting storage, Persistent Volume is bound to a Persistent Volume claim. Persistent volumes are 1:1 mapped to claims.

→ PVC takes the help of Storage class to create a P.V.

In case of Azure, storage class can be one of these 2

PV will be created and attached to your pod.

Now, CCM will create a 815GB of disk depending on which storage class you have asked for.

Now, when writing a Pod (remember, we used host=emptydir) → here you will

be using something called as

- name: volume  
persistentVolumeClaim:

claimName: azure-managed-disk.

volume mounts:

- mountPath: "/mnt/azure"  
name: volume.

and  
we can mount it into  
wherever we want

## Storage ~~class~~ classes in AWS :-

- We have a storage class called GPF (general purpose 2)

↓  
it is a disk that gets created when you  
create an EC2 instance.

→ `storage → mkdir Persistentvolumes & cd:-`

1) mysql-PVC.yaml

\$ kubectl get PV

2) mysql-pod.yaml

O/p: No resources.

\$ kubectl api-resources | grep storageclass

\$ kubectl get PVC

O/p:- storageclasses  
false

O/p: NO

SC  
↓  
short name.

\$ kubectl get SC

\$ kubectl api-resources | grep persistent

O/p: - - There are storage classes  
- - which are available.

O/p:- persistentvolumeclaims  
true

PVC v1

persistentvolumes  
false

PV v1

{ For me to create a claim, all I need  
is storage class.

\* Right now, the provisioning which we are using is Dynamic Provisioning → why?

I'm not creating storage. I'm directly writing PVC → there I'm making a request saying that I need 1 GB disk and then it is getting created.

\* Lets create a my-sql Pod with dynamic volume Provisioning :-

• Get storage classes. I got it with kubectl get sc

I can use any of the storage classes.

Google: az aks dynamic volume provisioning.

Let's Create a Persistent volume claim :-

PVC → is making a request to create volume dynamically.

API-ref → (version 1.75).



PersistentVolumeClaim v1 Core.

version = v1

kind = PersistentVolumeClaim.

metadata

↓  
object meta → name = mysql-PVC → This name is important → as we will be using this name in actual.

Spec

pe...v.claimSpec.

→ accessModes

It'll be using ReadWriteOnce



bcz I cannot have 2 mysql dbs trying to write at the same time bcz → it will mess.

(No disk will handle concurrency for you)

In K8S, whenever we are trying to P.Vs there are diff access modes.

ReadWriteOnce → this volume can be mounted as read-write

ReadOnlyMany

ReadWriteMany

by a single node.  
This access still <sup>can</sup> allows multiple pods to access the volume when the pods are running on the same node.

The volume can be mounted as read-only by many nodes.

The volume can be mounted as read-write by many nodes

For Eg: If you put your Excel file in shared disk and then if you open, — (85)  
If someone else tries to open → it will try to tell that you cannot write anything to this file bcoz there's some lock that has been created. That lock is not created by the disk. It is created by the Excel app.

But if you open notepad app? → it will not complain → whether 2 people open or 200 people open from the shared disk bcoz → disk cannot do anything. It is <sup>done by</sup> that app?. For that we are writing access mode.

⑥ **ReadWriteOnce** → we are trying to tell → now whenever you are mounting this → mount this to the node and ensure that any pod who <sup>are basically running</sup> ~~is~~ on the same node can mount it as a volume.

**RWOP** → **ReadWriteOnePod** is also available.

↓  
means give only to one pod & not to node.

→ resources → requests:

→ storageClassName: ?                   storage: 1Gi                   for static Provisioning  
   ↓  
   managed                           Google it.

### \* mysql-prc.yaml :-

---

apiVersion: v1

kind: PersistentVolumeClaim

metadata:

name: mysql-prc

spec:

accessModes:

- ReadWriteOnce

resources:

  requests:

    storage: 1Gi

storageClassName: managed

## \* mysql-pod.yaml :-

---

apiVersion: v1

kind: Pod

metadata:

: name: mysql-dynamic-vol

: labels:

: app: mysql

: layer: db

spec:

containers:

- name: mysql

: image: mysql:8

: ports:

- containerPort: 3306

: volumeMounts:

- name: mysql-prol

: mountPath: /var/lib/mysql

: volumes:

- name: mysql-prol

: persistentVolumeClaim:

: claimName: mysql-pvc

}

will automatically create a disk for me.

> kubectl apply -f \mysql-pvc.yaml

> kubectl get pvc

O/p: mysql-pvc status  
Pending.

> kubectl get pv

O/p: No resources

> kubectl apply -f \mysql-pod.yaml

O/p: Pod/mysql-dynamic-vol created.

> Kubectl get pvc

<u>Name</u>	<u>Status</u>	<u>Volume</u>	<u>Capacity</u>	<u>Access modes</u>	<u>Storage class</u>
mysql-pvc	Bound	pvc-1a...cebc	1Gi	RWO	managed

Is it bound now? = YES.

Until and unless you are using that - it will not create a disk.

> Kubectl get pv

It created 1 GiB disk now.

<u>Name</u>	<u>Claim</u>	<u>Capacity</u>	<u>Access Mode</u>	<u>Reclaim Policy</u>	<u>Status</u>	<u>Storage class</u>
pvc-1a...cebc	default/mysql-pvc	1Gi	RWO	Delete	bound	managed

1st I have created a PVC.

In default namespace,  
I have mysql-pvc.

After that, I asked for pvc → it said - it is in Pending state

↓

Then, I have applied my-sql-pod

↓

Once I applied mysql-pod → if I asked - it said it is bound.

\* what we have done is not very effective one. We'll be changing that.

We just created some storage space. Let's try to put it to action.

> Kubectl ~~delete~~ -f .

> Kubectl get pvc

O/p: No resources

> Kubectl get pv

O/p: pvc-1a...cebc 1Gi RWO Status Storage  
released managed.

I still have my disk. I can still use it.

\* Now, I want to create mysql pod. While creating it, it has to set some environmental variables → to try to tell what is mysql user, mysql <sup>root</sup> password etc.

\* Now, let's modify the spec → to add mysql related environment variables.

Exercise:

Try to create any Pod spec by setting Environment Variable.

env is an array of EnvVars.

↳ ↳  
name  
value.

\* mysql-pod.yaml :-

---

apiVersion: v1

kind: Pod

metadata:

name: mysql-dynamic-vol

labels:

app: mysql

layer: db

Spec:

containers:

- name: mysql

image: mysql:8

ports:

- containerPort: 3306

volumeMounts:

- name: mysql-pvol

mountPath: /var/lib/mysql

env:

- name: MYSQL\_ROOT\_PASSWORD

value: rootroot

- name: MYSQL\_USER

value: qtdevops

- name: MYSQL\_PASSWORD

value: rootroot

- name: MYSQL\_DATABASE

value: employees

volumes:

- name: mysql-pvol

persistentVolumeClaim:

claimName: mysql-pvc

(E1: Alpine Pod with  
'Env.Var → 'PURPOSE = TEST'  
with cmd → sleep 1d.)

writing passwords in plain text is  
not a good idea

(we will solve it later)

> kubectl apply -f .

\$ kubectl exec -it mysql-dynamic-vol -- mysql -uroot -p

O/p: Enter Password: rootroot

= = I'm inside mysql

mysql> show databases;

O/p:

Database
employees
information-schema
mysql
Performance-schema
sys

Google  
mysql sample create table Example  
w3schools.

mysql > use employees;

O/p: Database changed.

mysql insert into.

mysql > CREATE TABLE Persons (

-> PersonID int,  
-> LastName varchar(255),  
-> FirstName varchar(255),  
-> Address varchar(255),  
-> City varchar(255)  
-> );

mysql > INSERT into Persons (PersonID, LastName, FirstName, Address, City) VALUES (1,

'one', 'one', 'one', 'one')

(2,

'one1', 'one1', 'one1', 'one1')

mysql > Select \* from Persons

O/p:

PersonID	LastName	FirstName	Address	City
1	One	One	One	One
2	One1	One1	One1	One1

mysql > exit

Now, delete Pod and recreate.

\$ kubectl delete pod mysql-dynamic-vol

> kubectl apply -f mysql-pod.yaml (wait for pod to be created & running).

> kubectl exec -it mysql-dynamic-vol -- mysql -uroot -P

O/p: Enter Password: rootroot

mysql> use employees;

Database changed.

mysql> Select \* from Persons;

O/p:-

PersonID	LastName	FirstName	Address	City
1	One	One	One	One
2	One1	One1	One1	One1

The data is still available even after we deleted and recreated the Pod.

So, can I persist the data? → Yes.

→ Go to Azure → myResourceGroup → myAKS cluster → we can see PVC.

↓  
disk that got created (1 GB)

• Actually I have written as PVC

↓

that has spoken with StorageClass → and actually created

(Azure is something which has ) → (since there is a CCM a disk in Azure.  
(nothing to do with what K8s is) → (→ it is creating the stuff for you).

\* Writing a Pod might not be a good idea. Because Pod cannot maintain its state.

So, we need some controller.

One controller which we can use for this is → we can try to write replicaset and give one replica.

But there is some problem with that approach.

We have a better solution for this.

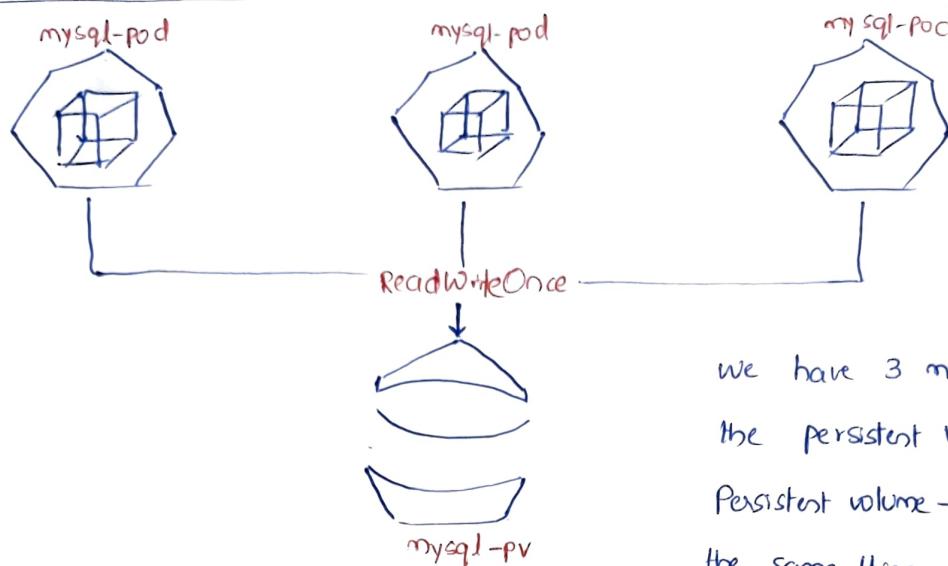
> kubectl delete -f .

-> If you want to scale pods we can use \* deployment  
\* replicaset.

\* You will have 3 pods . You will have a mysql-pv and all these pods

Our basic architecture :-

are using  
mysql-pv.



so, it is better to create these

kind of things with ReadWriteOnce.

means give Read Write Permission  
to only one node.

I can also give

one more access

= ReadMany → I don't have a problem with multiple pods reading.

But write Permission → I want to give it only to one pod.

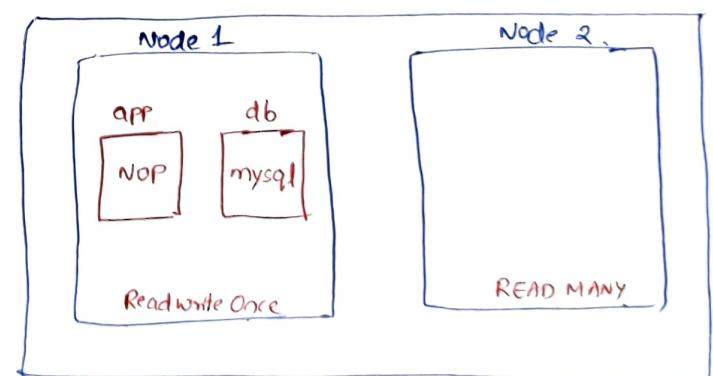
How you might have difficulties figuring out → which Pod has write permission & which Pod has read permission.

To avoid this problem

People don't give Readmany.

we have 3 mysql pods which use the persistent volume.

Persistent volume → if they are sharing the same thing → that means they have the same data . But there will be concurrency problems over here.



## Stateful Sets

In stateful sets, basically what would happen is you would have a Pod.

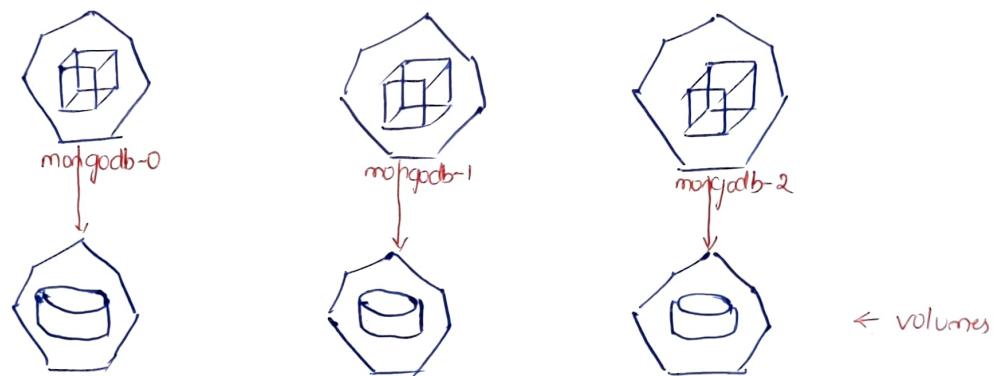
Generally, if you try to look at replica sets → In replicaset, whenever it creates a replica — it gives a random name to the pods (E.g.: mysql-some digits)

↳ mysql-0      }  
                mysql-1      }  
                mysql-2      }  
replicas=3  
} → It'll be getting  
procedural names (or)  
predictable names.

(Now, whenever, you try to create a P.V → you will be getting a volume for every Pod. → in Stateful sets)

- Stateful Set is like deployment with replicas. But each pod gets its own volume.
- Stateful Set is for stateful applications.

### Stateful Set



Each pod gets its own persistent volume. This is something which we call as a Database cluster.

So, we have basically an app which is

running in Pod which has its own state — which is basically maintained.

Now, for some reason, if you delete mongo-0 pod, then → a new Pod will be created which exactly gets mounted to the same volume.

Here each pod gets its own volume. They don't rely on 3 pods using same volume.

- If you remember, K8s stores everything in etcd cluster. All the modern dbs are designed to be distributed. This is a distributed system.

\* Now, here what happens is → whenever you try to use the data, data gets stored in its own storage (Pod's):

### Elastic Search :-

In Elastic Search, we have a concept called **shards**.

So generally what happens is → whenever you store the data in these kind of databases → what they do is → they don't store the entire data at one place. They break that data into multiple pieces and they store in multiple things.

For Ex:-

lets assume, we have data - 'A'. This 'A' data gets broken into number of shards  
 ↓  
 is a record

shards is how do you want to break this record.

So, I want to break it into 5 equal pieces.

A

↓

A1, A2, A3 A4 A5.

and

You can configure how many replicas you want.

Replica A1 = RA1.

④ You can configure → saying that → for every document that got broken → I want to have 5 shards and 2 replicas.

So, whatever you are storing is broken into 5 pieces and after replication → you are playing with 15 pieces.

Let's assume → 9 hark 3 databases → A, B, C.

A-node		B-node		C-node	
A1	A2	A3	RA1	RA2	RA3
A4	A5	RA1			
RA2	RA3	RA4			
RA5	RA6	RA7			
RA8	RA9	RA10			

→ A3 gets stored here.

→ RA1 gets stored here.

Let us assume that, if your B-node is failed

↓  
 what are the parts which you don't have.  
 I can get it from other node.

So, at any moment, even if one node goes for a toss,  
 you ~~can~~ have your whole doc. retrieved  
 like this.

→ The process of breaking your one record into multiple smaller portions is called as **Shards**.

For every shard → You will create replicas.

So, if you are having a 3-node cluster — If you try to create 5 shards and 2 replicas → You will not lose data even if 2 nodes fail.

This is how modern databases work.

• So for us to understand this functionality,

I'll be trying to create a stateful set of nginx.

\* So, when we create replicas in statefulset we get Predictable names.

\* We can access individual pod, by creating headless service and by using ..... svc.cluster.local.

\* Now I can individually access everything and then I'll be able to get the info. of that pod.  
bcz I have headless service in front of it.

We all know that headless service does not have an IP address. ~~Then~~ It gives the IP addresses of the pods. and we are using <Pod-name>..... svc.cluster.local.

\* Let us experiment with stateful sets and create **nginx pods** :-

We know that nginx is stateless. but I'll be trying to do something special in nginx.  
whatever we do now → we can apply it to any database.

→ Storage → <sup>new folder</sup> statefulsets → nginx-statefulset.yaml ;—

\$ kubectl api-resources | grep statefulset

Qp: true      sts      apps/v1

→ Docker Hub → nginx

**docker playground**  
\$ docker container run -d --name nginx -P nginx  
\$ docker container -it nginx /bin/bash.

Port: 49153 → → 0 → welcome to nginx ✓

90

Root

Let us see from which folder its coming from

root: # ls /usr/  
Op: bin/ lib/ share/

# cd /usr/share/nginx/html

# ls  
Op: index.html  
# cat index.html  
Op:- <!DOCTYPE html>  
<html>  
<head>  
<title> Welcome to nginx !</title>  
<style> ....</style>  
</head>  
<body>  
<h1>Welcome to nginx !</h1>  
...  
</body>  
</html>

Over here lets try to create → test.html

# vi test.html

Op: command not found.

# echo "<h1> i found it </h1>" > test.html

Now go to home page → ip:32-: /test.html → i found it.

I'll be trying to create 3 nginx-pods and then ill be creating 3 html files with 3 different values.

Then let us try to delete nginx pods by using kubectl delete commands.  
It will recreate the pod and let us see whether this exists (or) not.

• So, I'm trying to create a volume which mounts this → /usr/share/nginx/html.

Ideally, what we are about to do → should be used for **stateful application**.

(nginx became a stateful app bcs I'm storing that html page)

which is specific to this pod and then for that I'm trying to create a volume.

In reality, you will not be using this for nginx. You'll be using it for databases.

→ API reference → StatefulSet v1 apps



apiVersion: apps/v1

kind: StatefulSet

metadata

↳ name = nginx.

spec



type: StatefulSetSpec.

↳ minReadySeconds = 1

↳ PersistentVolumeClaimRetentionPolicy → describes the lifecycle of PVCS created from

volumeClaimTemplates.

↳ PodManagement Policy

↳ controls how pods are created  
during initial scale up.

By default, all PVCS are created as needed and  
retained until manually deleted.  
↳ we want the same.

The default policy is 'OrderedReady' where pods are created in increasing order

(pod-0, then pod-1 etc) and the controller will wait until each pod is ready  
before continuing. When scaling down, the pods are removed in opp order.

The alternative policy is 'Parallel' which will create pods in parallel  
to match the desired scale without waiting and on scale down will

delete all pods at once.

Default is  
ok for us.

replicas = 3

selector

LabelSelector.

matchLabels:

app: nginx

template

metadata  
name → nginx

labels:  
app: nginx

spec:

containers  
- name: nginx  
image: "  
ports... 80

volumeClaimTemplates:

- metadata:

spec:

storageClassName: managed

- - -

## nginx-svc.yaml :-

apiVersion: v1

Kind: Service

metadata:

name: nginx-svc

spec:

selector:

app: nginx

type: ClusterIP

clusterIP: None

ports:

- name: nginx-svc

port: 80

targetPort: 80

protocol: TCP

I need to give service name in statefulset.yaml.

NOW if you want to access any Pod

↓  
Pod name - Service name .

> kubectl get po

O/p: alpine 1/1 . . . → Let alpine pod be running → bcoz we have work with that pod.

> kubectl apply -f .

O/p: service/nginx-svc created

statefulset.apps/nginx created.

> kubectl get svc

O/p: clusterIP: none

> kubectl get sts -w

O/p: nginx ready

0/3

:  
3/3

It is taking time  
(bcoz PV needs to be created  
and it needs to be mounted.)

so we are creating disk also.

If it is a stateless app → we don't need to do all of this.

Ideally you should be running databases over here.

## nginx-statefulset.yaml :

```

apiVersion: apps/v1
kind: StatefulSet
metadata:
  name: nginx
spec:
  minReadySeconds: 1
  serviceName: nginx-svc
  replicas: 3
  selector:
    matchLabels:
      app: 'nginx'
  template:
    metadata:
      name: nginx
      labels:
        app: 'nginx'
    spec:
      containers:
        - name: nginx
          image: nginx
      ports:
        - containerPort: 80
      volumeMounts:
        - name: htmlhome
          mountPath: /usr/share/nginx/html
  volumeClaimTemplates:

```

## template

### volumeClaimTemplates:

#### - metadata:

name: htmlhome

#### spec:

##### accessModes:

- ReadWriteOnce

##### resources:

##### requests:

storage: 1Gi

storageClassName: managed

This will create a  
Persistent volume claim  
(PVC)  
Per pod.

Now, I'll get 3 pods with name

nginx-0, nginx-1, nginx-2.

/# curl nginx-0-nginx-svc

O/p: curl: not found

/# aptk add curl

/# curl nginx-0.nginx-svc

O/p: <html>

<head> <title> 403 Forbidden </title> </head>  
<body>

403 Forbidden.

— —

</html>

It says that it is forbidden.

/# nslookup nginx-svc

O/p: Name: nginx-svc.default.svc.cluster.local

Address: 10.244.1.16

Name: — — —

Address: 10.244.1.17

": — — —

: 10.244.0.40

( API-ref  
Look into service naming  
StatefulSet v1 core  
Pods get  
DNS/hostnames that follow the pattern:  
Pod-specificString.ServiceName.default.svc.  
cluster-local .

/# curl http://nginx-0.nginx-svc.default.svc.cluster.local

O/p: Forbidden.

↓  
cd /usr/share/-

ls.

→ There is no index.html. → That is the reason

we created an Empty volume.

/# curl http://nginx-0.nginx-svc/sts.html

O/p: <h1>Pod0-STS</h1> ✓

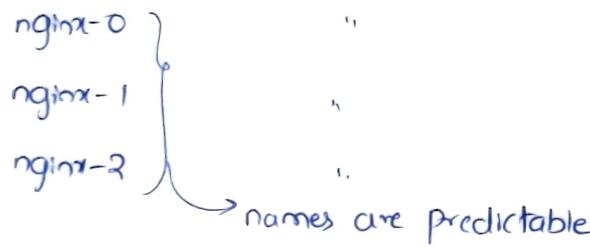
<Pod1> -1 ✓

<Pod2> 2 ✓

\* Now create a different html page for every pod.

\$ kubectl get po

O/p: alpine                    Running



Naming is Predictable  
in Stateful Sets

\$ kubectl delete po nginx-0

O/p: Pod 'nginx-0' deleted.

\$ kubectl get po

	Ready	Status
alpine	1/1	running
nginx-0	0/1	ContainerCreating
nginx-1	1/1	running
nginx-2	1/1	running

It is not creating nginx-0.  
It will again try to create nginx-0.

This is the major difference when it comes to deployment of stateful set

\$ kubectl exec nginx-0 --

Get into Pod & Create.

\$ kubectl exec -it nginx-0 -- /bin/bash

# cd /usr/share/nginx/html

# echo "<h1> Pod0-STS </h1>" > sts.html

# exit.

do same for nginx-1 & for nginx-2.  
Pod1-STS                    Pod2-STS

So, I have created 3 different html files in 3 different pods which are stored on 3 different persistent volumes.

\$ kubectl get svc

O/p: name -

\$ kubectl exec -it alpine -- /bin/sh

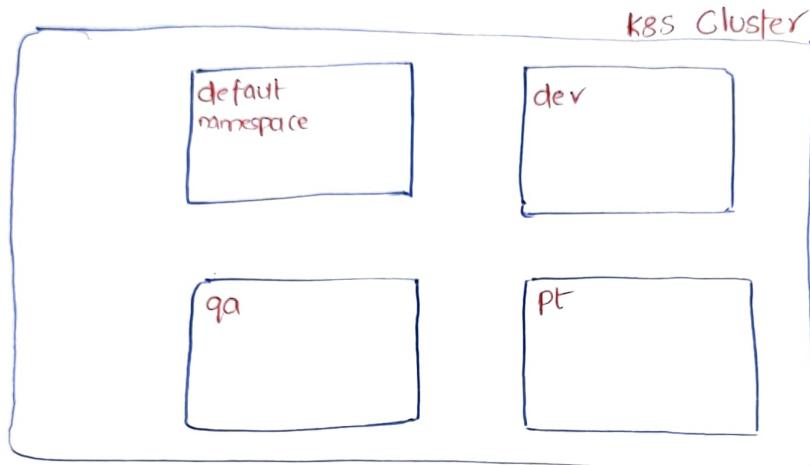
#

\* Now, some of the k8s resources will get created inside this namespace and some don't get created inside this namespace.

● In reality, what your Org<sup>zns</sup> will do is → They will try to create dev namespace

QA namespace, PT namespace

So this namespace is like a Virtual Cluster.



\$ kubectl api-resources

Eg: we can see there are some resources where `namespace=true` and some resources where `namespace=false`.

→ Namespace = true means → they are created within that namespace.

For Eg: pods → for pods `namespace` value is true.

But

For Eg: PersistentVolumeClaims → namespace is false.  
↓

what it means is → when you create a k8s pod → they belong to some namespace

As we saw just now → to access k8s service → we have used

service-name.default.svc.cluster.

↓  
name of your namespace.

Eg: Services → for services → `namespace=true` → so whenever you are creating a service → it exists in a namespace.

\* We sent a curl request from alpine pod.

> kubectl delete po nginx-0

O/p: Pod deleted

Wait for Pod to be up

> kubectl get po -w

O/p: ✓✓✓

Let's send request to nginx-0 from alpine.

# curl http://nginx-0.nginx-svc/sts.html

O/p: <h1>Pod0-STS</h1>

→ So, even if I deleted the pod-it was recreated and we got the same data back.

## #### Kubernetes namespace :-

• Generally whenever we are working with K8S, there are various namespaces

\$ kubectl get ns

O/p: default

• Any namespace that starts with 'Kube-'

gatekeeper-system

is something which is created by K8S

kube-node-lease

→ Gatekeeper-system → is the namespace created by Azure

kube-public

to do some of its activities.

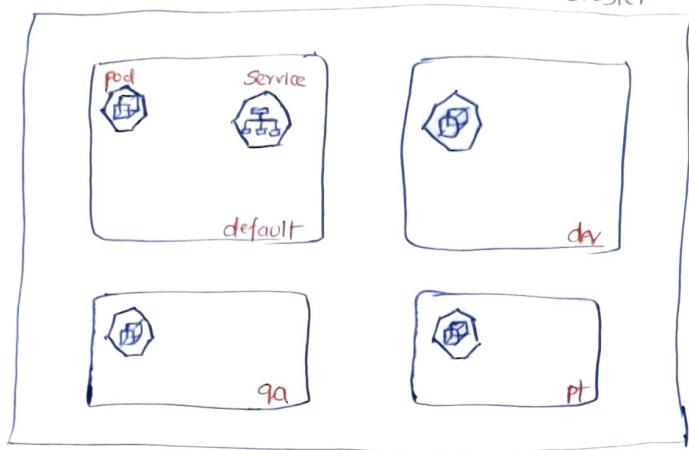
• All these days we were using a namespace called Default.

\* Let's assume, we have 3 nodes and K8S cluster. We are assuming that pods get created in these nodes in cluster.

But in reality it doesn't work that way.

In each of <sup>the</sup> K8S cluster, what K8S does → it tries to create something called as namespace. This what we're trying to speak about is a default namespace.

## K8S Cluster



So, if you want to access within this cluster →

you have to give → resource-name • namespace • svc.cluster.local.

$\downarrow$   
this is fully qualified domain name.

\* The basic idea is to separate one section of pods with other sections.

> kubectl delete -f .

> kubectl create ns dev

O/p: namespace/dev created.

> kubectl get po

O/p: alpine 1/1 Running → It'll be getting pods. But it'll be getting only pods of bydefault namespace which is 'default'.

> kubectl get po --namespace default

O/p: alpine 1/1 running (All these days whenever you were executing the cmd → it was equivalent to this.)

> kubectl get po --namespace kube-system

O/p:- Name                  Ready                  Status . . .

ama-logs-85984                  3/3

ama-logs-fw                  3/3

ama-logs-rs                  2/2

= =                  :

= -                  :

To run your K8S cluster, there are already lot of pods doing the job for you.

But these are all created by K8S cluster. we don't need to worry about it.

> kubectl get po --all-namespaces

O/p: — You have these many pods which are running across diff. namespaces.

• But, for the user, acc to the permissions which he has, he has permissions on particular namespaces.

\* You don't give any namespace → it will use 'default'.

\$ kubectl get po --namespace dev

O/p: No resources found.

In metadata, there is one interesting thing called as namespace.

(Namespace defines the space within which each name must be unique.)

(An empty namespace is equivalent to "default" namespace).

Eg:- nginx-statefulset.yaml → apiVersion: apps/v1  
nginx-svc.yaml  
apiVersion: v1  
kind: Service  
metadata:  
name: nginx-sv  
namespace: dev  
kind: StatefulSet  
metadata:  
name: nginx  
namespace: dev  
I, (in dev namespace)

(I'm creating this stateful set in a namespace called dev.)

> kubectl apply -f .

O/p: created

\$ kubectl get sts

O/p: No resources found in default namespace.

(Bcz it is querying only default namespace)

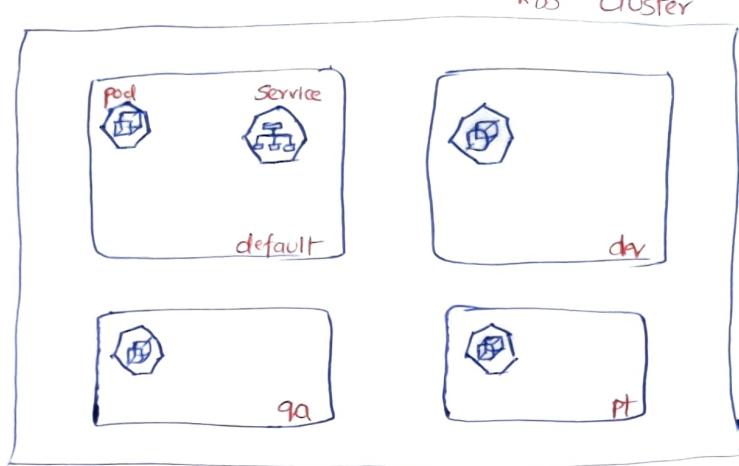
\$ kubectl get sts --namespace dev

O/p: 

Name	Ready
nginx	1/3 19s

\* So, in your org<sup>nz</sup>, basically what wud happen is you will have a production k8s cluster, staging(k8s) cluster and then you wud have one k8s cluster for dev, test- and all testing environments.

\* Your org<sup>nz</sup> ask you to specify the namespace. Depending on which envt you want to deploy - you have to pass that namespace.



So, if you want to access within this cluster →

you have to give → resource-name • namespace • svc.cluster.local.

\* This ↓ is fully qualified domain name.

The basic idea is to separate one section of pods with other sections.

> kubectl delete -f .

> kubectl create ns dev

O/p: namespace/dev created.

> kubectl get po

O/p: alpine 1/1 running → I'll be getting pods. But I'll be getting only pods of by default namespace which is 'default'.

> kubectl get po --namespace default

O/p: alpine 1/1 running (All these days whenever you were executing the cmd → it was equivalent to this.)

> kubectl get po --namespace kube-system

O/p:-	Name	Ready	Status	- - -
	ama-logs-85984	3/3	7/7	To run your k8s cluster, there are already
	ama-logs-fw	3/3	7/7	lot of pods doing the job for you.
	ama-logs-rs	2/2	7/7	But these are all created by k8s cluster.
	= =	:	7/7	We don't need to worry about it.
	= -	:	7/7	

> kubectl get po --all-namespaces

O/p: --- You have these many pods which are running across diff. namespaces.

\$ kubectl get po

O/p: alpine

\$ kubectl exec -it alpine -- /bin/sh

I want to access the other pod. How do I do that?

/# curl nginx-0.nginx-svc.dev.svc.cluster.local

O/p: test.



Now lets directly try to access this

/# curl nginx-0.dev.cluster.local

O/p could not resolve host

/# nslookup nginx-0.dev  
Server: 10.0.0.10, Address: 10.0.0.10:53

/# nslookup 127.0.0.1

O/p: Server can't find 1.0.0.127... : servfail.

Full name includes



pod name . namespace name . cluster-local

since I'm using headers sent over here



( )

)

Google: pods fqdn

↓  
DNS for services and pods.

④ pod-ip-address.my.namespace-pod.cluster-domain.example. ( In general a pod has the )

if a pod is in default namespace.

fol. DNS resolution

Ex:- if domain name for your cluster is = cluster.local.

then pod has a DNS name



172-17-0-3.default.pod.cluster.local.

→ Any pods exposed by a Service have the foll. DNS resolution available.

↓  
Pod-ip-address.service-name.my-namespace.svc.cluster-domain.example.

Google: get fqdn in linux.

↓  
hostname --fqdn.

/# hostname --fqdn

O/p: nginx-0.nginx-svc.dev.svc.cluster.local

( Since there is a headless service over here → It gives headless service's name )  
but actually → You'll get ip-addr stuff

> kubectl get ~~ps~~<sup>sts</sup> --namespace dev

O/p: nginx 2/3 56s

But →

\$ kubectl get pv

- Persistent volumes do not belong to any namespace.

So, they are shared resources.

- So, whenever you are seeing  $\rightarrow$  with namespaced value as false.

\$ kubectl api-resources →

↓  
They don't belong to any particular

That means that

They are not different acc. to the namespaces.

namespace.

- So, Persistent Volume Claim can be created within the namespace, but persistent volume can be shared across namespaces.

So, for dev namespace as well as for test namespace  $\rightarrow$  I might be using same persistent volume.

$\rightarrow$  So, this namespace is an abstract idea  $\rightarrow$  where you can create clusters within cluster and this is a logical cluster.

Note: Not every K8S resource will be part of namespace.

For e.g.: volumes are not part of namespace.

- So, if you want to communicate from a pod which is running in one namespace to the other namespace  $\rightarrow$  You have to use the fully qualified domain name.

i.e.  $\cdot <\text{pod-name}> \cdot <\text{namespace-name}> \cdot \text{svc} \cdot \text{cluster.local}$ .

> kubectl exec --namespace dev -it nginx-0 -- /bin/sh

# cd /usr/share/nginx/html

# ls

o/p: lastfound. There are nothing in this.

# echo "Test" > index.html

# ls .-1-> index.html

# \*ETCD Database backup and restore\*

→ ETCD :-

- \* ETCD is a distributed reliable key-value store that is simple, Secure & fast.
- \* The ETCD database stores info. regarding the cluster such as nodes, pods, configs, secrets, accounts, roles, bindings etc.
- \* Every Info. you see when you run the 'kubectl get' command is from the ETCD database.
- \* ETCD starts a service that listens on port 2379 by default.

# etcdctl -version

O/P: 3.3.13

API version: 2.

We need API version 3 for this.

# kubectl -n kube-system get pod

O/P: 

Name	Ready	Status
coredns		

etcd-controlplane running 1/1

{ kube-apiserver-controlplane

  kube-controller-manager-controlplane

  kube-flannel-ds

  kube-proxy

  kube-scheduler-controlplane

→ All are deployed as static pods (done by kubeadm tools)

# kubectl -n kube-system describe pod etcd-controlplane

We get all the info. regarding key, certificates

Volumes

etcd-certs: Path: /etc/kubernetes/pki/etcd  
HostPathType: DirectoryOrCreate

etcd-data: Type: HostPath

Path: /var/lib/etcd

HostPathType: DirectoryOrCreate

# kubectl get deployments.apps

O/p: blue 3/3  
red 2/2

# kubectl -n kube-system describe pod etcd-controlplane

O/p:-

etcd:

-- cert-file = /etc/kubernetes/pki/etcd/server.crt ✓

-- data-dir = /var/lib/etcd .

-- name = controlplane

-- peer-cert-file = /etc

-- key-file = /etc/kubernetes/pki/etcd/server-key ✓ ✓

-- trusted-ca-file = /etc/kubernetes/pki/etcd/ca.crt . ✓

Google :-

Kubernetes

↓  
search: etcd snapshot

Command:

ETCDCTL\_API=3 etcdctl --endpoints=https://127.0.0.1:2379 \

--cacert=<trusted-ca-file> --cert=<cert-file> --key=<key-file> \

snapshot save <backup-file-location>

Eg: /opt/snapshot-pre-boot.db

# ETCDCTL\_API=3 etcdctl --endpoints=https://127.0.0.1:2379 \

> --cacert=/etc/kubernetes/pki/etcd/ca.crt --cert=/etc/kubernetes/pki/etcd/server.crt  
--key=/etc/kubernetes/pki/etcd/server.key \

> snapshot save /opt/snapshot-pre-boot.db

( Backup of etcd is now taken )

O/p:- Snapshot saved at /opt/snapshot-pre-boot.db

NOW something disastrous has happened and etcd is not responding

↓  
# kubectl get deployments.apps

O/p: No resources found → data lost. How to recover ?

We can just restore etcd.

Command for restoring ETCD :-

ETCDCTL\_API=3 etcdctl --endpoints 10.7.0.9:2379 snapshot restore   
we are moving our db to particular locations.  
So - we will use  
--data-dir=/var/lib/etcd-backup.

 ↓  
/opt/snapshot-pre-bootdb

# ETCDCTL\_API=3 etcdctl --data-dir=/var/lib/etcd-backup snapshot restore /opt/snapshot-pre-bootdb

O/P:- etcdserver/membership: added member 8e9...94d [http://localhost:2380] to cluster cd.c32

# cd /var/lib

# ls

O/P: apt config contained docker dockerdocker etcd  git kubelet . . .

↓ we can see that folder has been created.

# cd /etc/kubernetes/manifests/

Here we have all config files (static pods)

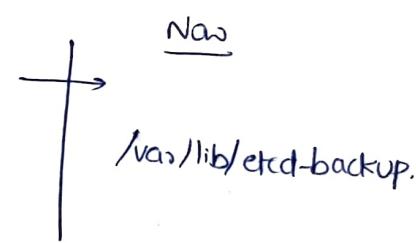
# ls

O/P: etcd.yaml kube-apiserver.yaml kube-controller-manager.yaml kube-scheduler.yaml

# vi etcd.yaml.

You need to change the path →

Previous  
volumes:  
path: /var/lib/etcd



When we save this file

it will do proper linking.

# docker ps | grep etcd

There is nothing pending — linking has been established.

# kubectl get deployments.apps

O/P: blue 3/3 . . .

red 2/2 . . .

# \* Kubernetes Cluster Upgrade \*

- At a time Kubernetes supports only 3 minor version.
- Currently, Main version is 1 and minor version is 21. (1-21) On average they release 4 versions within a year
- K8s releases minor versions very frequently in a year.
- We will be upgrading the cluster using Kubeadm utility.
- Upgrade will be carried out one node at a time, starting with master node.
- We need to drain the master/worker node before doing upgrade.

## → Master node Upgrade Command :-

\* Kubeadm upgrade plan

\* Kubectl drain controlplane -ignore-daemonsets

\* apt update

\* apt install kubeadm=1.19.0-00

\* kubeadm upgrade apply v1.19.0

\* apt install kubelet=1.19.0-00

\* systemctl restart kubelet

\* Kubectl uncordon controlplane. (making master-node ready to work)

## → Worker node Upgrade Command :-

→ Kubectl drain node01 -ignore-daemonsets

→ ssh node01

→ apt install kubeadm=1.19.0-00

→ Kubeadm upgrade node

→ apt install kubelet=1.19.0-00

→ systemctl restart kubelet

→ Kubectl uncordon node01

## Kubernetes upgrade from v1.18 to v1.19 :-

Lets monitor the Env first.

\$ kubectl get node

O/p: controlplane Ready <sup>version</sup>  
node01 not Ready v.1.18.0

Lets drain the Master-node first.

\$ kubectl drain controlplane --ignore-daemonsets

O/p: node/controlplane cordoned.

evicting pod kube-system/kube-dns

- - -

node/controlplane evicted

\$ kubectl get node

O/p: controlplane Ready, SchedulingDisabled  
node01 Ready.

\$ kubectl get pod → we will find only pods from node

O/p: goldenginx-8d-running:-node1 & not find any pod from Control plane.

Lets start the Upgrade process.

\$ apt update

\$ apt install kubeadm=1.19.0-00

\$ kubeadm upgrade apply v1.19.0

O/p: Now that your control plane is upgraded, please proceed with upgrading your kubectls if you haven't already done so.

\$ apt install kubelet=1.19.0-00

\$ systemctl restart kubelet

\$ kubectl get node

O/p: controlplane v1.19.0 Ready, schedulingdisabled  
node01 v1.18.0 Ready

Now uncordon to allow master node to work.

\$ kubectl uncordon controlplane

O/p: node/controlplane uncordoned.

\$ kubectl get node

	status	version
controlplane	Ready	v1.19.0
node01	Ready	v1.18.0

\$ kubectl drain node01 --ignore-daemonsets

Now its evicting all the pods we have on node01. To control plane component

\$ kubectl get pod -o wide

Name	Ready	Status	Node	Nominated node
gold-ngrx...	1/1	running	control-plane	<none>

Pod which was running previously one node is now running on controlplane.

so, we moved it, bcoz we want to upgrade worker node.

controlplane

\$ kubectl get node

O/p:	controlplane	Ready	v.1.19.0
	node01	ready, schedulingDisabled	v.1.18.0

\$ ssh node01

node1 \$ apt update

\$ apt install kubeadm=1.19.0-00

\$ kubeadm upgrade node

\$ apt install kubelet=1.19.0-00

\$ systemctl restart kubelet

ctrl+d to come out from Worker node.

controlplane \$ kubectl get node

O/p:	controlplane	Ready	v.1.19.0
	node01	Ready, schedulingDisabled	v.1.19.0

\$ kubectl uncordon node01

\$ kubectl get node

O/p: Both are Ready. ✓