

Date: \_\_\_\_\_

# Kubernetes (K8s)

↳ **Fireship**

managing containerized workloads in cloud

e.g. Kubernetes for binance - market closed

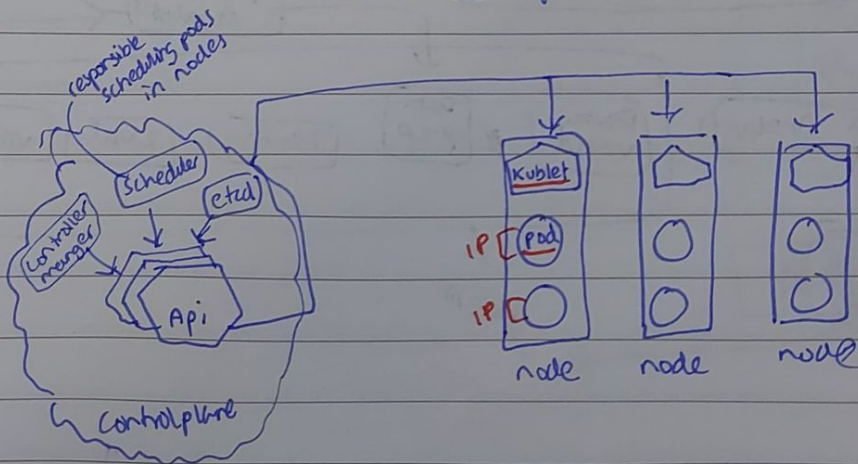
- market open → scale containers  
replace container if one full

A system deployed on Kubernetes is called a **Cluster**

Brain of operation is the **control plane** which exposes a **api server**

↓  
also contains it Key/Value database  
called **ETCD** - store info about running  
the cluster or persistent state

which can handle internal  
external req to manage cluster



→ Control plane manages worker machines called **Nodes**

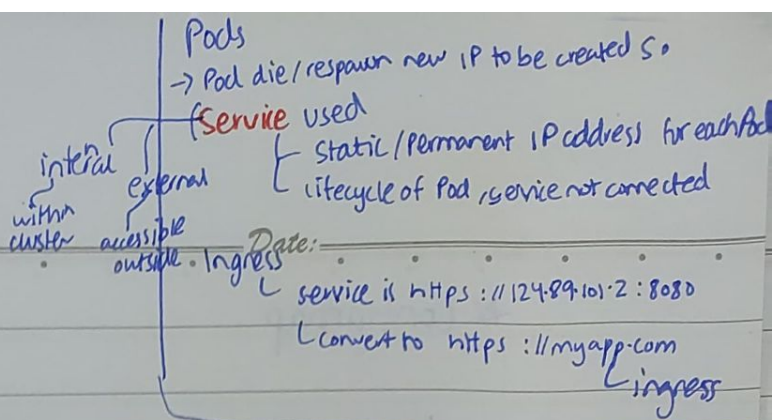
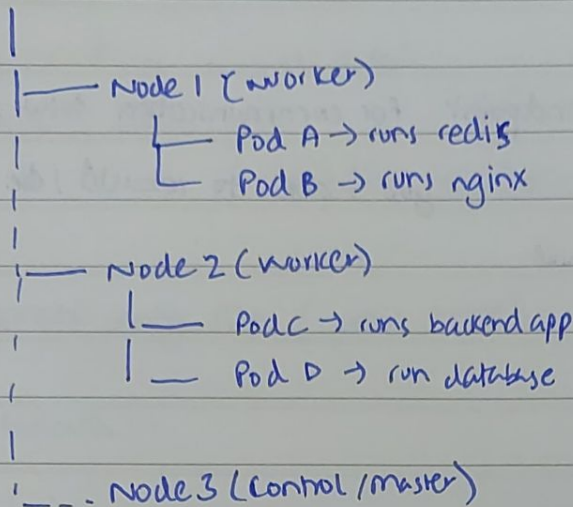
→ each Node runs a **Kublet** - tiny application to communicate back with the main control plane mothership

→ each node contains multiple **Pods** - smallest deployable unit in K8s  
↳ group of containers / 1 or more

**Node** | **Pod**  
- machine where containers run | - 1 or more containers

provides share storage / network for containers  
pods ephemeral = dies span new IP  
each pod has own IP address

e.g. K8 Cluster



- workload increases then K8 increase <sup>no of</sup> ~~more~~ nodes 2
  - also takes care of
    - networking
    - secret management
    - persistent storage ...
- K8 is designed for High Availability
  - achieves by maintaining
- You describe **objects** in yaml
  - define state of cluster

replica set:

set of running pods ready to go at any time

```

apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx-deployment
spec:
  replicas: 3
  containers:
    - name: nginx
      image: nginx:1.14.2
      ports:
        - containerPort: 80
      volumes:
        - name: cool-volume
  
```

// describe behavior in spec field

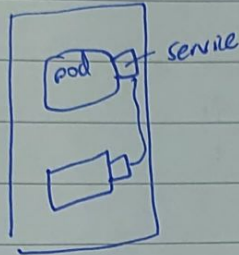
replica controller



Date: .....

## ★ Configmap

if pod changed the endpoint for communication b/w pods change and thus you have to rebuild the image, make pod



solution: configmap

↳ outside of image just rename new pod's name with endpoint

DB-url = mongo-db-service  
          endpoint          pod-name

## ★ Secrets

↳ alike configmap but it user, password of a database

newpod = diff user, password so use secret to store these

stored in base64

use 3rd party encryption for further security

## ★ Volume

↳ ~~same~~ stores persistent data

↳ outside of pod  
references ~~out~~  
the storage.

↳ local  
↳ remotely

★

Deployment — replicas of pod  
↳ abstraction above pods  
↳ if pod dies substitute with replica

deployment config file. (prev page)

```
{  
  metadata: {  
    spec: {  
      status: // running status etc.  
    }  
  }  
}
```

- Problem: we can't replicate ~~DB~~ Db pod to ensure data consistency  
so use stateful set



Date: \_\_\_\_\_

# DVC data versioning control

git can't store large dataset, model in ML

so we use dvc (tool) to track this

dvc ~~tracks~~ tracks real dataset/model  
git tracks metadata files of dataset/model

• git + dvc used in collaboration

↳ links to an external drive (Gdrive, S3)

Base project

```
ml-project/  
├── data/  
│   └── train.csv  
├── model/  
│   └── r  
├── src/  
│   ├── preprocessing.py  
│   └── train.py  
├── params.yaml  
├── readme.md  
└── runmodel.py
```

## Example workflow

1) git init // makes .git folder  
dvc init // makes .dvc folder

# Track dataset

2) dvc add data/train.csv // track dataset

• **metadata** file appears → train.csv=dvc

↳ data/  
└── train.csv  
    └── train.csv.dvc

• Data copied to .dvc/cache

git add data/train.csv.dvc • gitignore  
git commit

✓  
dvc config core.autostage true // automatically track file changes  
↳ auto ~~adds~~ runs dvc add

## \* Tracking changes

1. change data : a script like get-data.py changes data

2. check status :

{ git doesn't find changes  
but dvc does dvc status

3. Add change to dvc `dvc add data.csv`

4. Git with note `data.csv.dvc` has changed add it

5. commit

## \* Go back to previous data

`dvc logs` // find prev version

`git checkout`

`dvc checkout`

## \* Remote storage setup

Add remote: `dvc remote add g-drive-remote gdrive://(folder-url)`  
This command changes config file so git add, commit  
do get client-ID, secret from google

To push: `dvc push -r gdrive-remote`  
`dvc pull`



Date: \_\_\_\_\_

file automate  
getdata, preprocess, train

## Automation: Dvc pipeline, dvc.yml

dvc → used to automate pipeline to retrain model - when new data come

dvc repro // reruns pipeline

dvc metrics show

dvc params show

dvc metrics diff // show differences in metrics

dvc.yml eg

stages:

get-data:

cmd: python src/get-data.py

deps:

- src/get-data.py

outs:

- data/raw.csv

preprocess:

...

train: ...

How it knows to rerun

dvc checks

1- scripts (deps)

2- datafile (outs)

3- parameters (params.yml)

for changes

if found

dvc repro

## Automate commits or CI/CD using Github Actions

1- Run dvc repro

2- Add updated files

3- commit, push to

create .github/workflows/auto-train.yml file

page

Date:

name: Autot train pipeline

on:

push:

branches:

- main

pull request:

jobs:

run-pipeline:

run-on: ubuntu-latest

steps:

- name: Set-up python

uses: actions/setup-python

with:

python-version: 3.10

- name: Install dependencies

run: pip install requirement.txt

- name: Pull data

run: dvc pull --force

- name: rerun pipeline

run: dvc repro



Date: . . . . .

# MLflow

↳ tool/library in python to manage ml experiments

helps in

- Tracks experiments
- Reproduce training runs anywhere
- Save / Load models easily
- version & manage models

→ 1  
→ performance score (e.g accuracy)  
→ files (plots, save models)

★ Components	Purpose	e.g
Tracking	log experiments (parameter, metrics, artifacts)	compare model / tuning results
Projects	Package code so can run anywhere	Run training code anywhere
Models	Save models in standard format	deploy model anywhere
Model Registry	version manage models	promote best model to production

## 1. Tracking

import mlflow  
import mlflow, sklearn

with mlflow.start\_run():

mlflow.log\_param("learning\_rate", 0.01)

" ◦ log - metric

" ◦ log - artifact

mlflow.sklearn.log\_model(model)

\$ mlflow ui // start mlflow ui, go to localhost:5000

\$ mlflow autolog() // autodetects model & logs every param

### Workflow

1. install mlflow
2. import mlflow
3. initialize (set server IP)
3. train model as normal
4. start run
  - ↳ log

Tip

## hyperparameter tuning using gridsearch

↳ method to auto test multiple combinations of parameters to find best

↳ you give - model  
- dictionary of parameters

it returns  
best model + best parameters

gridsearch (model, param dictionary, cv = 5...)

Date: . . . . .

## 2, Projects

requires .git repo

This file

create MLproject file

```
name: my-ml-project
conda-env: conda.yml
entry point:
```

main:

parameters:

```
n_estimators: {type: int, default: 100}
max_depth: {
commented: "python train.py --n-estimators {n_estimators} . . . . ."
```

function that can be run

optional: conda.yml file for dependency

## \* Running the project

\$ mlflow run . // default param

\$ mlflow run . -P n\_estimators=200 -P max\_depth=8 //

## 3, models

when you save a model using mlflow it creates a repo like this

```
model/
├── mlmodel // metadata of model
├── model.pkl // actual model
└── conda.yml // dependencies
```

### \* model flavors

pytorch	mlflow.sklearn
scikit-learn	• pytorch
XG-boost	• "
Tensorflow	• "

example import mlflow.sklearn  
with mlflow.start\_run()

model = RandomForest

model.sklearn.log\_model(model)

mlflow saves model in standard format so

- can be
  - loaded into python
  - used in webapp
  - deployed as REST API

mlflow models serve -m runs:/{run-id}/model -p 1234

serve as API, now send request at localhost:1234

saved in

mlruns/

<experiment-id>/

<run-id>/

artifacts/

model/

└ model



## 4, Model Registry

↳ Hub of models

- ↳ track model version (v1, v2, ...)
- ↳ Promotes model : Staging → Production
- ↳ stores notes, metadata, description

```
result = mlflow.register_model("runs:/<run-id>/model", "HousePriceModel")
```

```
mlflow.transition_model_version_stage(
```

```
    name = "HouseP..."
```

```
    version = 1
```

```
    stage = "Production"
```

```
)
```

~~Demos on MLflow~~

- ran
- 1- python tracking.py
  - 2- mlflow run -e tracking --env-manager local

→ didn't run and  
as requires more info

→ we defined experiment name in tracking.py but didn't tell same name in cmd  
so it created another experiment & conflicted with H name

→ Add "--experiment-name Basic-Tracking-Demo" at end of  
cmd 2