# Contents

# Understanding the Kubernetes Architecture

# Cluster Architecture

SCHEDULER

POD → ASSIGN POD TO NODE

NODE NODE NODE

ETCD

DATA

CLUSTER CONFIG

MASTER

CONTROLLER MANAGER

MAINTAIN CLUSTER

NODE FAILURE

REPLICATE COMPONENTS

API SERVER

COMMUNICATION HUB

# API SERVER

## KUBELET
Manages containers
on the node

## KUBE-PROXY

### API SERVER

NODES [KP] [KP] [KP] [KP]

Network Proxy that runs on each node

Network rules on nodes

## WORKER NODE

## CONTAINER RUNTIME
→ DOCKER
→ CONTAINERD

RUNS YOUR
CONTAINERS

# Application running on Kubernetes

MASTER

NODE 1
- POD
- KUBELET
- KUBE-PROXY
- DOCKER

2 N

IMAGE REGISTRY

# API Primitives

API server is the only one that communicates with Etcd
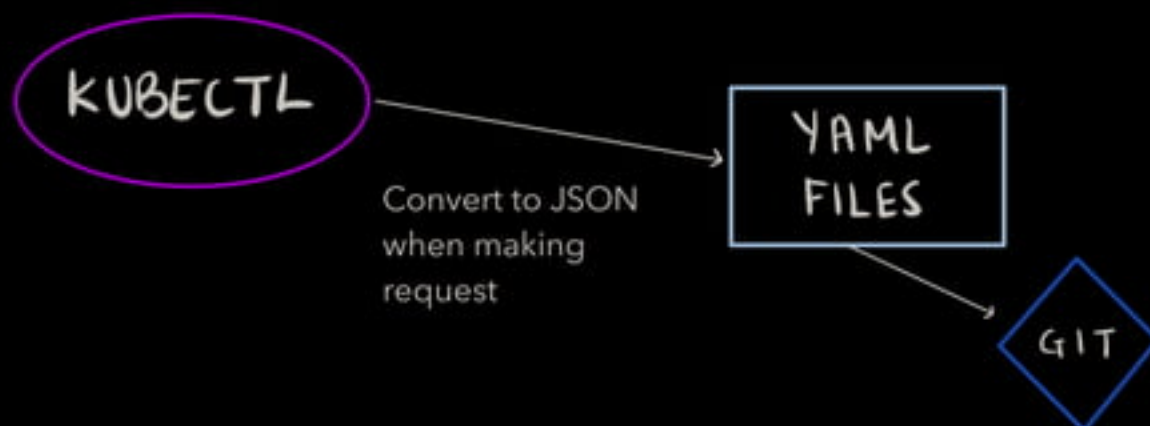
Every component communicates with the API server only and not directly with one another.

Objects like pods and services are declarative intents.

KUBECTL

Convert to JSON when making request

YAML FILES

GIT

# YAML File Composition

**API VERSION** ———————→ Clear consistent view of resources

**KIND** ———————→ The kind of object you want to create
- Pod
- Deployment
- Job

**METADATA** ———————→ Uniquely identify the object

Name String      UID      Namespace

**SPEC** ———————→ Container image volume exposed ports

**STATUS** ———————→ State of the object –> match desired states

# Services and Network Primitives

Services allow you to dynamically access a group of replica pods



NODE 1

NODE 2

32767

32767

SERVICE

POD

POD

POD

CONSISTENT IP ADDRESS

# Kube-Proxy

Kube-Proxy handles the traffic associated with a service by creating IP table rules

# Building
# the
# Kubernetes Cluster

# Release Binaries, Provisioning and Types of Clusters

PICKING THE RIGHT SOLUTION

CLOUD          OR          ON-PREM

## CUSTOM                              PRE-BUILT

- Install manually
- Configure your own network fabric        OR
- Locate the release binaries
- Build your own Images
- Secure cluster comms

- Minikube
- Minishift
- Micro K8S
- Ubuntu on LXD
- AWS, Azure and GCP

# Installing kubernetes master and nodes

## MASTER + WORKERS

① DOCKER +
KUBERNETES ————→ GPG KEY

————→ ADD REPOS

② UPDATE PACKAGES

③ INSTALL DOCKER, KUBELET, KUBEADM, KUBECTL

④ MODIFY BRIDGE ADAPTER SETTINGS

## MASTER ONLY

① INITIALISE CLUSTER

② MAKE DIRECTORY FOR K8S

③ COPY KUBE CONFIG

④ CHANGE OWERSHIP OF CONFIG

⑤ APPLY FLANNEL CNI

# Building highly available cluster

All components can be replicated, but only certain ones can operate simultaneously

# Replicating etcd

**TOPOLOGY**

## STACKED

EACH CONTROL PLANE NODE

CREATES LOCAL

**ETCD MEMBER**

ONLY COMMUNICATES WITH

**API SERVER**

## EXTERNAL TO KUBERNETES CLUSTER

ETCD — ETCD → **RAFT CONSENSUS ALGORITHM**

ETCD

**API SERVER**

**REQUIRES MAJORITY**

THERE MUST BE MORE THAN HALF
TAKING PLACE IN THE STATE
CHANGE AND THEREFORE THERE
MUST BE ODD NUMBER OF NODES

# Configuring Secure Cluster Communications

→ ALL COMMUNICATION VIA HTTPS

→ KUBECTL —TRANSLATES→ KUBERNETES API

PROVIDES CRUD → CLUSTER STATE

CRUD:
- CREATE
- READ
- UPDATE
- DELETE

RETURN RESPONSE

KUBECTL CREATE POD —HTTP POST→ API SERVER → STATE STORED IN ETCD

MULTIPLE PLUGINS

[ AUTHENTICATION | AUTHORISATION | ADMISSION | VALIDATION ]

READ → SKIP

CALLS TO DETERMINE REQUEST
- HTTP HEADER
- CERTIFICATE

CAN THIS USER PERFORM THIS ACTION?

CREATE MODIFY DELETE

# Building Highly Available Cluster

RBAC is used to prevent unauthorised users from modifying the cluster state

ROLE BINDING → ROLE → ACTION

WHO CAN DO IT

1 OR MORE

WHAT CAN BE DONE

RESOURCE

| ROLE | ROLE BINDING | → NAMESPACE RESOURCE |

| CLUSTER ROLE | CLUSTER ROLE BINDING |

CLUSTER LEVEL RESOURCES

## Service Account

POD ← SERVICE ACCOUNT ← IDENTITY OF APP RUNNING IN POD → API

/VAR/RUN/SECRET/ KUBE·IO/SERVICE ACCOUNT

TOKEN FILE ← AUTHENTICATION TOKEN

**NAMESPACE 1**

POD   POD

SERVICE ACCOUNT

**NAMESPACE 2**

POD   POD

SERVICE ACCOUNT

ONLY USE SERVICE ACCOUNT IN SAME NAMESPACE

# Running end to end tests on cluster

PERFORMANCE
AND
RESPONSE OF
APPLICATION

← WHY? →

EXAMPLE
TESTS

KUBE TEST

POOR CLUSTER
PERFORMANCE

- Deployments can run
- Pods can run
- Pods can be directly accessed
- Logs can be collected
- Commands run from pods
- Services can provide access
- Nodes are healthy
- Pods are healthy

# Managing Cluster

UPGRADING CLUSTER → KUBEADM

→ KUBECTL

OPERATING SYSTEM UPGRADES

USE DEPLOYMENTS OR REPLICA SETS

GET PODS
DRAIN NODE
UNCORDON
↑
PUT BACK TO SERVICE

BACKUP AND RESTORE CLUSTER → CLUSTER STATE → ETCD → SAVE EXTERNALLY

# Network
# Cluster
# Communication

# Pod and Node Networking

## NETWORKING WITHIN A NODE

10.244.1.2        10.244.1.3

POD 1       POD 2

ETH0         ETH

VETH2AA      VETH8D8

BRIDGE

10.244.1.1/24

NODE 1

## NETWORKING OUTSIDE OF THE NODE

10.244.1.2              10.244.2.3

POD 1 → VETH       VETH ← POD 2

BRIDGE → ETH0     ETH0 ← BRIDGE

NODE 1 → 172.31.43.91     NODE 2 → 172.31.34.149

NETWORK

# Container Network Interface

10·244·1·2

( POD 1 ) → VETH
          ↓
    BRIDGE → ETH0 — ( CNI ) ← ETH0 ← BRIDGE
                                        ↑
NODE 1 → 172·31·43·91           VETH ← ( POD 2 )

10·244·2·3

NODE 2 → 172·31·34·149

**NETWORK**

CNI IS A NETWORK OVERLAY
  └→ ALLOWS BUILDING TUNNEL BETWEEN NODES
        └→ SITS ON TOP OF EXISTING NETWORKS
              └→ ENCAPSULATES PACKETS   | HEADER | DATA | TRAILER |
                    └→ CHANGES SOURCE/DEST   └— ADD —┘

HOW DOES
CNI DO
THIS?        ——→ THERE IS A MAPPING ASSOCIATED IN USER SPACE
                    └→ PROGRAM ALL POD IP ADDRESS'S TO NODE IP,
                       WHEN REACH OTHER NODE, DE-ENCAPSULATE
                       PACKET AND GIVE TO BRIDGE

EXAMPLE CNI

                              ⭐  APPEARS LOCAL TO NODE!

CALICO    FLANNEL

# Service Networking



**API SERVER**
SERVICE
10·109·185·62

**NODE 1** → 172·31·43·91
10·244·1·2
POD 1 → VETH
KUBE·PROXY
BRIDGE — IP TABLES — ETH0

**NODE 2** → 172·31·34·149
10·244·2·3
KUBE·PROXY
VETH ← POD 2
ETH0 — IP TABLES — BRIDGE

NETWORK

→ PODS COME AND GO
HOW DOES THE CLUSTER KEEP TRACK?
↳ SERVICES
↳ PROVIDES VIRTUAL INTERFACE → AUTO ASSIGNED TO PODS BEHIND INTERFACE

EXAMPLE

CLUSTER IP SERVICE → AUTO CREATED ON CLUSTER CREATION
↳ TAKES CARE OF INTERNAL ROUTING
↳ NO MATTER WHERE MOVES, OTHER PODS KNOW HOW TO COMMUNICATE TO IT

# Ingress Rules and Load Balancers

LOAD BALANCER

- REDIRECT TRAFFIC TO ALL NODES AND PORTS
- CLIENTS TALK TO LB TO ACCESS APPLICATION

SERVICE → TALKS TO 1 NODE AT A TIME
→ CANNOT SPLIT TRAFFIC LIKE A LOAD BALANCER

ONLY ACCESSIBLE INTERNALLY → DOES NOT HAVE EXTERNAL IP

LOAD BALANCER ← 1 EXTERNAL IP ADDRESS FOR EVERY SERVICE

NODE 1
POD 1
31732

NODE 2
POD 2
31732

SERVICE

# Ingress

INGRESS
→ APP·EXAMPLE·COM/APP 1 → SERVICE  POD POD
→ APP·EXAMPLE·COM/APP 2 → SERVICE  POD POD
→ WEB·EXAMPLE·COM → SERVICE  POD POD

ACCESS MULTIPLE SERVICES WITH SINGLE IP ADDRESS

## Cluster DNS

EVERY SERVICE DEFINED IN THE CLUSTER IS ASSIGNED A DNS NAME

SERVICE NAME        NAMESPACE        BASE DOMAIN NAME

JENKINS · DEFAULT · SVC · CLUSTER · LOCAL

50 · 10 · 0 · 1 · DEFAULT · POD · CLUSTER · LOCAL

POD IP        NAMESPACE        BASE DOMAIN NAME

A PODS DNS SEARCH WILL INCLUDE THE PODS OWN
NAMESPACE AND THE CLUSTERS DEFAULT DOMAIN

# Pod Scheduling

# Configuring the Kubernetes Scheduler

Scheduler responsible for assigning pod to node based on resource requirements of the pod

RULES ARE PLACED BY DEFAULT → HOWEVER, CAN CREATE OWN

WHY?

SAME NODE TO SAVE COSTS

WORKER NODES HAVE DIFFERENT DISKS

SCHEDULER →

| | |
|---|---|
| 1 | DOES THE NODE HAVE ADEQUATE HARDWARE RESOURCES? |
| 2 | IS THE NODE RUNNING OUT OF RESOURCES? |
| 3 | DOES THE POD REQUEST A SPECIFIC NODE? |
| 4 | DOES THE NODE HAVE A MATCHING LABEL? |
| 5 | IF POD REQUEST A PORT, IS IT AVAILABLE? |
| 6 | IF POD REQUESTS A VOLUME, CAN IT BE MOUNTED? |
| 7 | DOES THE POD TOLERATE THE TAINTS OF THE NODE? |
| 8 | DOES THE POD SPECIFY NOD OR POD AFFINITY? |

# Running multiple schedulers for multiple Pods

It is possible to have 2 schedulers working alongside each other.

```
    POD    POD              POD    POD
         |                       |
         v                       v
  +-------------+         +-------------+
  | SCHEDULER 1 |         | SCHEDULER 2 |
  +-------------+         +-------------+
         |                    |      |
         v                    v      v
NODES -> ( 1 )            ( 2 )  ( 3 )
```

# Scheduling pods with limits and label selectors

TAINTS ———→ REPEL WORK ———→ EXAMPLE
                                MASTER NODE
                                NO SCHEDUAL

TOLERATIONS ——→ ALLOW YOU TO ——→ EXAMPLE ←—— DAEMON SET POD
                 TOLERATE A       KUBE-PROXY    MUST RUN ON ALL
                 TAINT                          NODES

(CPU/MEMORY) ——→ POD MAY NOT BE          SCHEDULER
                 USING ALL REQUESTED         |
                 RESOURCE AT A GIVEN          v
                 TIME                     POST 8 STEPS
     |                                    /          \
     v                          MOST REQUESTED    LEAST REQUESTED
SCHEDULER LOOKS                 PRIORITY          PRIORITY
AT THE SUM OF                        ↑
RESOURCES REQUESTED                  |            CLOUD ENVIRONMENTS
BY EXISTING PODS              WHY? ←——————————— YOU ARE PAYING FOR
                                                 ALL RESOURCES

# DaemonSets

DaemonSets ensure that a single replica of a pod is running on each node at all times



POD  DAEMON SET POD

POD  REPLICA SET POD

IF YOU TRY DELETE A DAEMONSET POD, IT WILL SIMPLY RECREATE IT

# Display Scheduler Events



LOG LEVEL ← SCHEDULER PROBLEMS → POD LEVEL

EVENT LEVEL

# Deploying Applications

# Deploying an Application, Rolling Updates and Rollbacks

DEPLOYMENTS ⟶ | HIGH LEVEL RESOURCE FOR DEPLOYING AND UPDATING APPS |



**DEPLOYMENT**

- REPLICASET
  - POD
  - POD
  - POD — APP: V1
- REPLICASET
  - POD
  - POD
  - POD — APP: V2

KUBECTL APPLY ⟶ MODIFY OBJECTS TO EXISTING YAML AND IF DEPLOYMENT NOT CREATED ⟶ ALSO CREATE

KUBECTL REPLACE ⟶ REPLACES OLD WITH NEW AND OBJECT MUST EXIST

ROLLING UPDATE ⟶ PREFERRED WAY ⟶ SERVICE NOT INTERRUPTED
⟶ FASTEST WAY

KUBECTL ROLLOUT ⟶ ROLLBACK PREVIOUS VERSION

# Configuring an App for HA and Scale

AVOIDING BAD DECISIONS → BLOCK BAD VERSION RELEASE

**MIN READY SECONDS**
↓
HOW LONG A NEWLY CREATED POD SHOULD BE READY BEFORE CONSIDERED AVAILABLE

AND

**READINESS PROBE**
↓
DETERMINES IF A SPECIFIC POD SHOULD RECIEVE CLIENT REQUEST OR NOT

---

DEPLOYMENT V1 → MIN READY SECONDS → READINESS PROBE → POD ERRORS AT 5 SECONDS → POD NOT RELEASED

MIN READY SECONDS ↓ 10

READINESS PROBE ↓ MUST RETURN SUCCESS → CHECK EVERY SECOND / PORT 80

POD NOT RELEASED ↓ ROLL BACK VERSION

---

PASSING CONFIGURATION OPTIONS TO APP
↓
ENVIRONMENT VARIABLES
↑
STORE IN CONFIG MAP
CREATE SECRET AND PASS TO EV
↑ MULTIPLE CONTAINERS CAN USE SAME
↑ JUST UPDATE, NO NEED TO REBUILD IMAGE

CONTAINER

/ETC/CONFIG → CONFIGMAP VOLUME: CONFIG →

/ETC/CERTS → SECRET VOLUME: CERTS →

POD

CONFIGMAP: APPCONFIG

| KEY A | VOL 1 |
|-------|-------|
| KEY B | VOL 2 |

SECRET: APPSECRET

| CERT | VOL |
|------|-----|
| KEY  | VOL |

# Creating a self-healing app

ReplicaSets ensure that identically configured pods are running at the desired replica count

RECOMMENDED

DEPLOYMENTS ← REPLICA SETS → LOSING NODE

MANAGES REPLICASETS

HAS NO IMPACT ON APP

PODS ARE UNIQUE

POD DIES

REPLACED WITH SAME HOST NAME AND CONFIG

STATEFULSETS

HEADLESS SERVICE

UNIQUE PODS

CERTAIN TRAFFIC TO GO TO EACH POD

VOLUME CLAIM

NEEDS OWN STORAGE AS IT IS UNIQUE

# Managing Data

# Persistent Volumes

PODS ARE EPHEMERAL $\rightarrow$ POD TERMINATED
$\downarrow$
STORAGE TERMINATED

STORAGE MUST BE INDEPENDENT $\longrightarrow$ IF POD MOVES
$\downarrow$
STORAGE FOLLOWS

STORAGE CLASSES $\rightarrow$

**PERSISTENT VOLUME**

PROVISIONING

STATIC

CLUSTER ADMIN
CREATES AND
AVAILABLE FOR
CONSUMPTION

DYNAMIC
$\rightarrow$ PVC MUST REQUEST
A STORAGE CLASS

RESOURCE IN
CLUSTER

## Volume Access Modes

By specifying an access mode with your PV, you allow the volume to be mounted to one or many nodes, as well as read by one or many

VOLUME CAN ONLY
BE MOUNTED USING
ONE ACCESS MODE
AT A TIME

MOUNT CAPABILITY OF
NODE NOT POD

### ACCESS MODES

READWRITE ONCE

ONLY 1 NODE CAN
MOUNT THE VOLUME
FOR READ AND WRITE

READWRITE MANY

MULTIPLE NODE
CAN MOUNT FOR
READ / WRITE

READONLY MANY

MULTIPLE NODE
CAN MOUNT VOLUME
FOR READING

# Persistent Volume Claims (PVC)

PVC allows the application developer to request storage for the application, without having to know underlying infrastructure.

DEV → PVC → PV

STAYS WITH PV

CLUSTER ADMIN → ACTUAL STORAGE

**POD**
- → VOLUMES
- → MOUNT PATH

**PVC**
- → 1 GI
- → RWO

**PV**
- RECLAIM POLICY ↓ RETAIN

BOUND

RETAIN DATA IN VOLUME

COULD ALSO BE RECYCLE OR DELETE

DELETE CONTENTS OF VOLUME

DELETE UNDERLYING STORAGE

# Storage Objects

Volumes that are already in use by a pod are protected against data loss. This means even if you delete a PVC, you can still access the volume from the pod.

STORAGE OBJECT IN USE PROTECTION → PVC CANNOT BE REMOVED PREMATURELY

② STILL BOUND

PVC
FINALIZERS
→ PVC PROTECTION

POD

PV

③ DELETE POD

PVC DELETED

① GETS DELETED

| PROVISIONER | PARAMETER | RECLAIM POLICY |
|---|---|---|
| ↓ | ↓ | ↓ |
| AWS-EBS | GP2 | RETAIN |

STORAGE CLASS

# Applications with Persistent Storage

EXAMPLE



STORAGE CLASS

METADATA
  ↳ FAST

PVC

STORAGE CLASS NAME
  ↳ FAST
→ 100Mi
→ READWRITEONCE

DEPLOYMENT

REPLICA → 1
IMAGE → KUBESERVE·V1
VOLUMEMOUNT → /DATA
VOLUME → VOLUME -DATA
PERSISTENTVOLUMECLAIM

→ ROLLOUT

PV

1. Create storage class object
2. Create PVC object
3. Create deployment
4. Rollout deployment
5. Check pods
6. Create file on mount
7. List contents

# Securing
# the
# Kubernetes Cluster

# Service accounts and users

API SERVER

FIRST EVALUATES

SERVICE ACCOUNT — I.E JENKINS

NORMAL USER
- PRIVATE KEY
- USER STORE
- FILE → USER/PASS LIST

KUBECTL CREATE SERVICEACCOUNT JENKINS →

**ASSIGN TO POD BY PUTTING IN POD MANIFEST** → IF YOU DO NOT USE SPECIFIC IT WILL USE DEFAULT

CREATES SERVICE ACCOUNT → KUBECTL GET SA
↓
WILL SHOW DEFAULT + JENKINS

SECRET
↓
HOLDS THE PUBLIC CA OF API SERVER + JWT TOKEN

YAML FOR SERVICE ACCOUNT
KUBECTL GET SA JENKINS -O YAML
↓
KUBECTL GET SECRET + NAML ← SHOW SECRET NAME

THIS IS WHAT REQUEST WILL BE USED TO AUTH WITH API SERVER

BUSYBOX-YAML

NAME: BUSYBOX
...
SPEC:
SERVICEACCOUNTNAME: JENKINS

JENKINS SERVER
↳ ADDS K8S CLI PLUGIN + TOKEN
↓
NOW CAN CONTROL PODS

# Service accounts and users

USER

👤

ACCESS CLUSTER REMOTELY → MASTER ↗ CREATE CLUSTER ROLE BINDING

ADNAN
↓
CERT

MASTER ↓ SET CREDENTIALS →
KUBECTL CONFIG ①
SET-CREDENTIALS
-- USERNAME = ADNAN
-- PASSWORD = PASSWORD

KUBECTL CONFIG
SET-CLUSTER KUBERNETES ②  ← SET CLUSTER
-- SERVER = HTTP://1.1.1.1
-- CERTIFICATE = CERT

KUBECTL CONFIG ③  ← SET USER
SET-CREDENTIALS

CONTEXT → CAN BE USED TO CONNECT
TO MULTIPLE CLUSTERS

KUBECTL CONFIG SET-CONTEXT
KUBERNETES --CLUSTER = KUBERNETES
--USER.... --NAMESPACE

④
USE CONTEXT

KUBECTL GET NODES ← AS PER NORMAL
⑤

# Cluster Authentication and Authorisation

AUTHENTICATION ⟶ AUTHORISATION

AUTHENTICATION ↓
FIRST STEP IN
RECIEVING REQUEST
↓
WHO? ⟨ POD?
       HUMAN?

AUTHORISATION ↓
WHAT CAN
THEY DO? ⟶ RBAC RULES
↓
4 RESOURCES
OR 2 GROUPS

WHO CAN
DO IT
↓

WHAT CAN BE
PERFORMED ON ⟵ ROLES AND
WHICH RESOURCE    CLUSTERROLES

ROLE BINDINGS
AND
CLUSTER ROLE BINDINGS

ROLE ⟵ ROLE BINDING
NAMESPACE 1

ROLE ⟵ ROLE BINDING
NAMESPACE 2

ROLE ⟵ ROLE BINDING
NAMESPACE 3

CLUSTER ROLE BINDING ⟶ CLUSTER ROLE

KUBERNETES CLUSTER

## ROLE

CREATE
NAMESPACE ⟶ CREATE ROLE ⟵ REFERENCE SINGLE ROLE ⟵ CREATE ROLE BINDING ⟶ CAN BIND TO MULTIPLE USER, SERVICE ACCOUNTS AND GROUPS

CREATE ROLE ↓
LIST SERVICES
FROM WEB ⟶ WHAT ACTIONS NOT WHO
NAMESPACE

## CLUSTER ROLE

CREATE
CLUSTER ⟶ CLUSTER ROLE BINDING
ROLE
↓
VIEW PERSISTANT VOLUMES ⟵
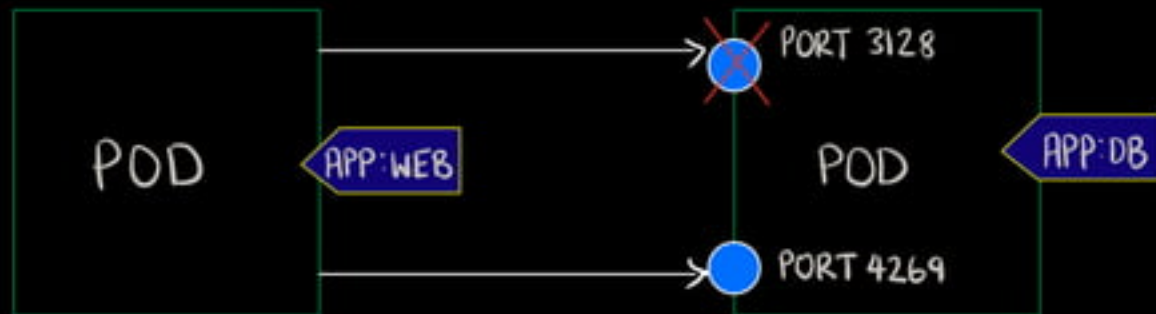
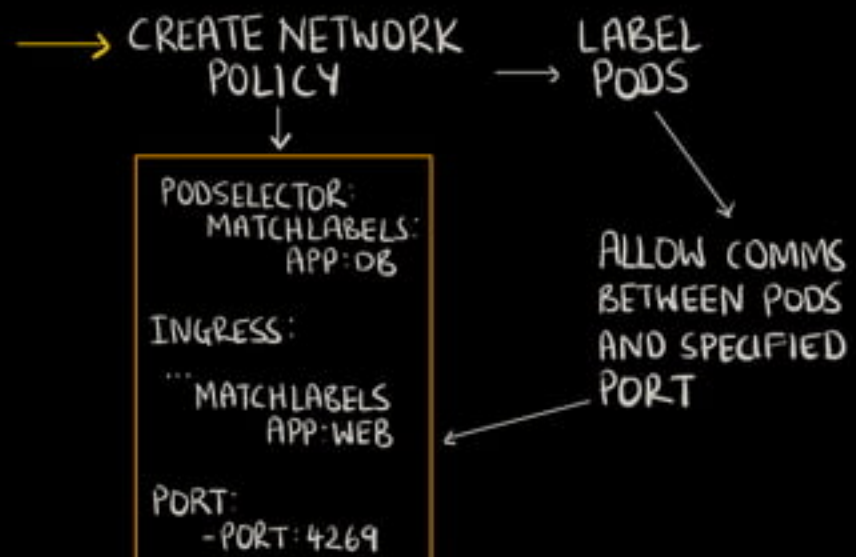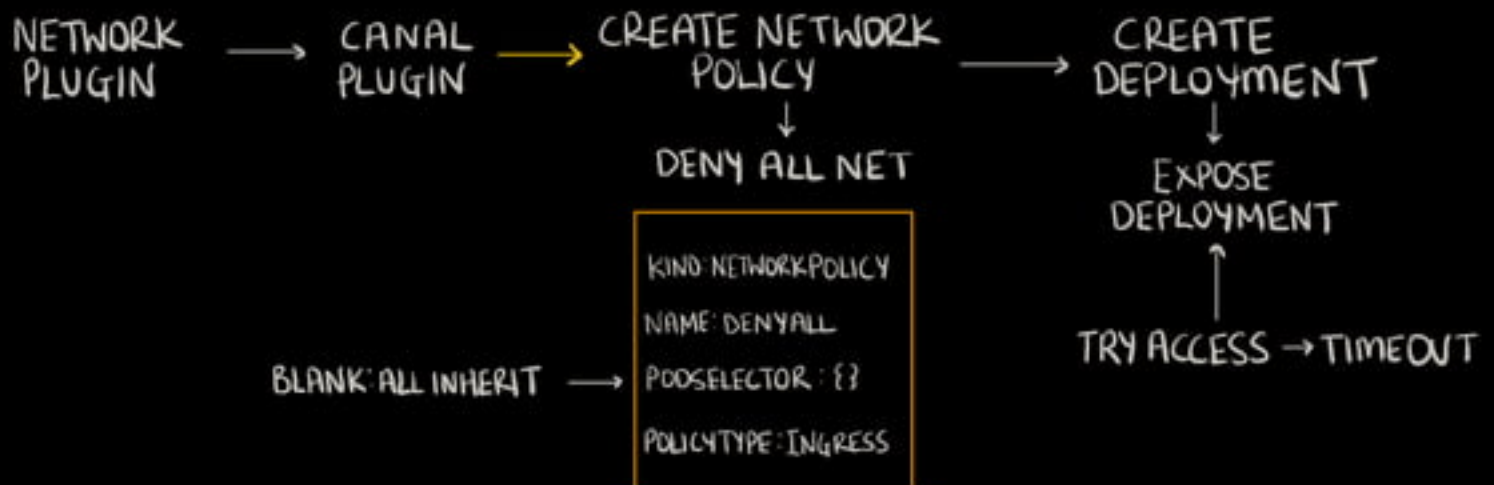POD
↓
CURL FROM ⟶ ACCESS AT
CONTAINER      CLUSTER LEVEL

# Configuring Network

Network policies use selectors to copy rules to pods for communication throughout the cluster

POD — APP:WEB → ✗ PORT 3128

POD — APP:DB → PORT 4269

## HOW?

NETWORK PLUGIN → CANAL PLUGIN → CREATE NETWORK POLICY → CREATE DEPLOYMENT

CREATE NETWORK POLICY ↓ DENY ALL NET

CREATE DEPLOYMENT ↓ EXPOSE DEPLOYMENT ↑ TRY ACCESS → TIMEOUT

BLANK: ALL INHERIT →
```
KIND: NETWORKPOLICY
NAME: DENYALL
PODSELECTOR : {}
POLICYTYPE: INGRESS
```

→ CREATE NETWORK POLICY → LABEL PODS

CREATE NETWORK POLICY ↓
```
PODSELECTOR:
    MATCHLABELS:
        APP: DB

INGRESS:
    ...
    MATCHLABELS
        APP: WEB

PORT:
    - PORT: 4269
```
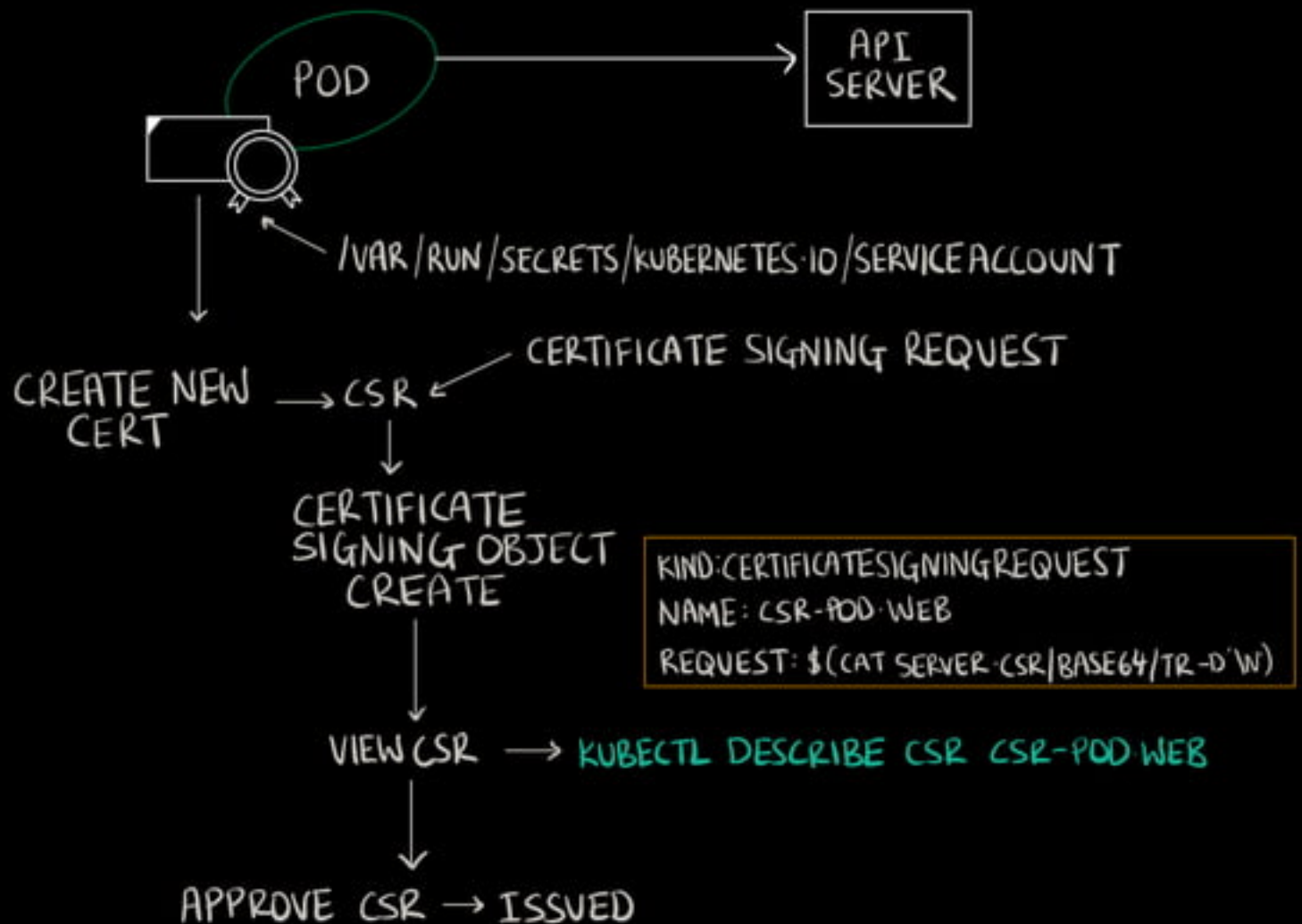
ALLOW COMMS BETWEEN PODS AND SPECIFIED PORT

# Creating TLS certificates

The CA is used to generate a TLS certificate and authenticate with the API server

POD → API SERVER

/VAR/RUN/SECRETS/KUBERNETES·IO/SERVICE ACCOUNT

CERTIFICATE SIGNING REQUEST

CREATE NEW CERT → CSR

CERTIFICATE SIGNING OBJECT CREATE

KIND: CERTIFICATESIGNINGREQUEST
NAME: CSR-POD·WEB
REQUEST: $(CAT SERVER·CSR|BASE64/TR-D'\W)

VIEW CSR → KUBECTL DESCRIBE CSR CSR-POD·WEB

APPROVE CSR → ISSUED

# Secure Images

<u>Private Registry</u>



CONTROLLING IMAGES THAT GO INTO PRODUCTION
→ VULNERABILITIES ⟶ CLAIR (SCANNING)
→ SOMETHING ON IT CAUSE NODE CRASH

→ LOGIN TO PRIVATE REGISTRY
　↳ TAG DOCKER IMAGE
　　↳ PUSH TO PRIVATE REGISTRY

## <u>HOW?</u>

USE FOR IMAGE PULLS

KUBERNETES CREATE SECRET → DOCKER REGISTRY TYPE → MODIFY SERVICE ACCOUNTS → POD

→ DOCKER-SERVER
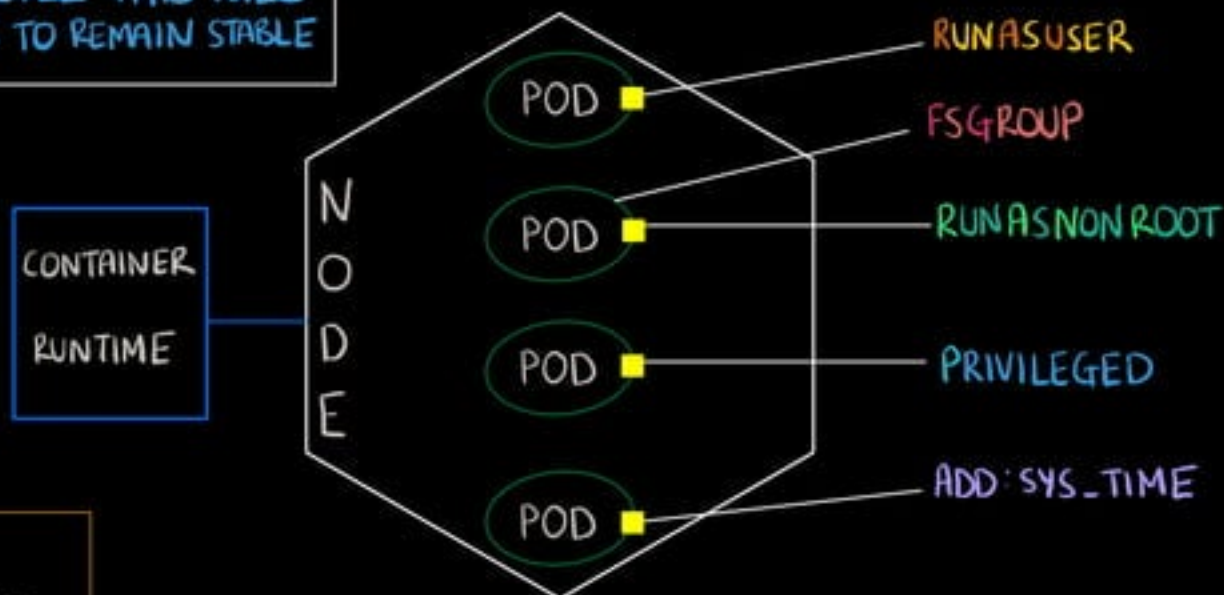→ DOCKER-USERNAME
→ DOCKER-PASSWORD

KUBECTL PATCH

POD ↓ PRIVATE REGISTRY IMAGE

# Defining Security Context

LIMIT ACCESS TO CERTAIN
OBJECTS AT THE POD AND
CONTAINER LEVEL THIS WILL
ALLOW IMAGES TO REMAIN STABLE

CONTAINER

RUNTIME

N O D E

POD → RUNASUSER

POD → FSGROUP

POD → RUNASNONROOT

POD → PRIVILEGED

POD → ADD: SYS_TIME

KIND: POD
IMAGE: ALPINE
SECURITYCONTEXT: → RUN POD AS 405
    RUNASUSER: 405

CAN ALSO PUT 'RUNASROOT'

ABILITY TO RUN AS PRIVILEGED → 'PRIVILEGED: TRUE'

## CONTAINER LEVEL

ABILITY TO LOCK
DOWN KERNEL
LEVEL FEATURES → SETTING
ON CONTAINER     CAPABILITIES
                 ON POD LEVEL

ADD → SECURITYCONTEXT:
        ADD:
          - SYS_TIME
          - NET_ADMIN

REMOVE → SECURITYCONTEXT:
           DROP:
             CHOWN

# Securing persistent key/value store

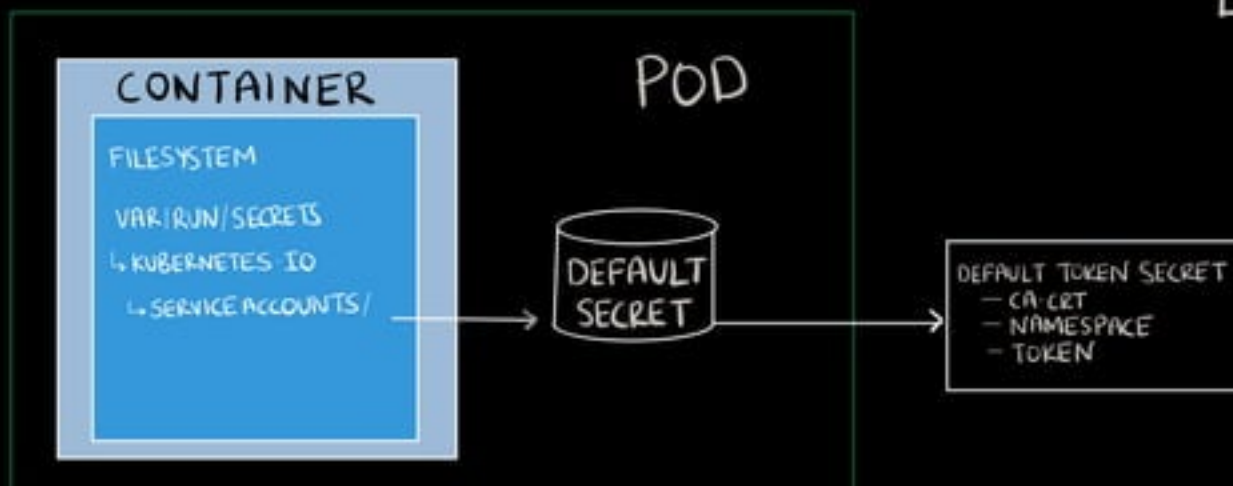Secrets allow you to expose entries as files in a volume keeping this data secure is crucial to cluster security

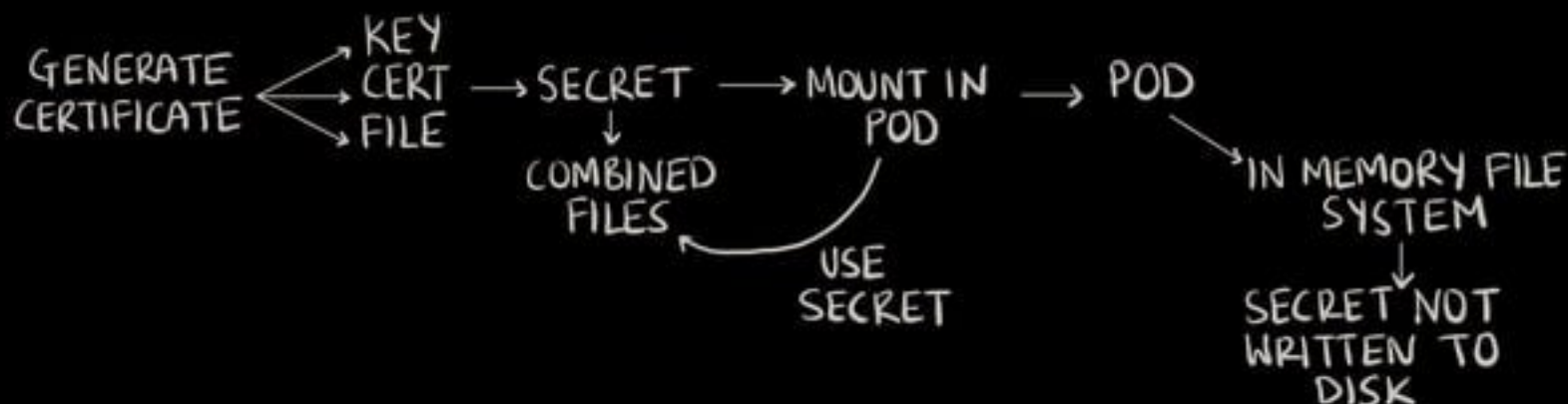DATA MUST LIVE BEYOND LIFE OF POD → SECRETS → KEY/VALVE PAIR

KEY/VALVE PAIR ↓ PASS AS ENV VAR ← NOT BEST PRACTICE
OR
EXPOSE AS FILES IN VOLUME

NOT BEST PRACTICE ↓ MAY BE OUTPUT TO LOG FILES

## POD

### CONTAINER

FILESYSTEM

VAR|RUN/SECRETS
↳ KUBERNETES IO
  ↳ SERVICE ACCOUNTS/

→ DEFAULT SECRET →

DEFAULT TOKEN SECRET
— CA·CRT
— NAMESPACE
— TOKEN

## HTTPS TO WEBSITE

GENERATE CERTIFICATE → KEY
CERT → SECRET → MOUNT IN POD → POD
FILE

SECRET ↓ COMBINED FILES

MOUNT IN POD ⤸ USE SECRET

POD ↘ IN MEMORY FILE SYSTEM ↓ SECRET NOT WRITTEN TO DISK
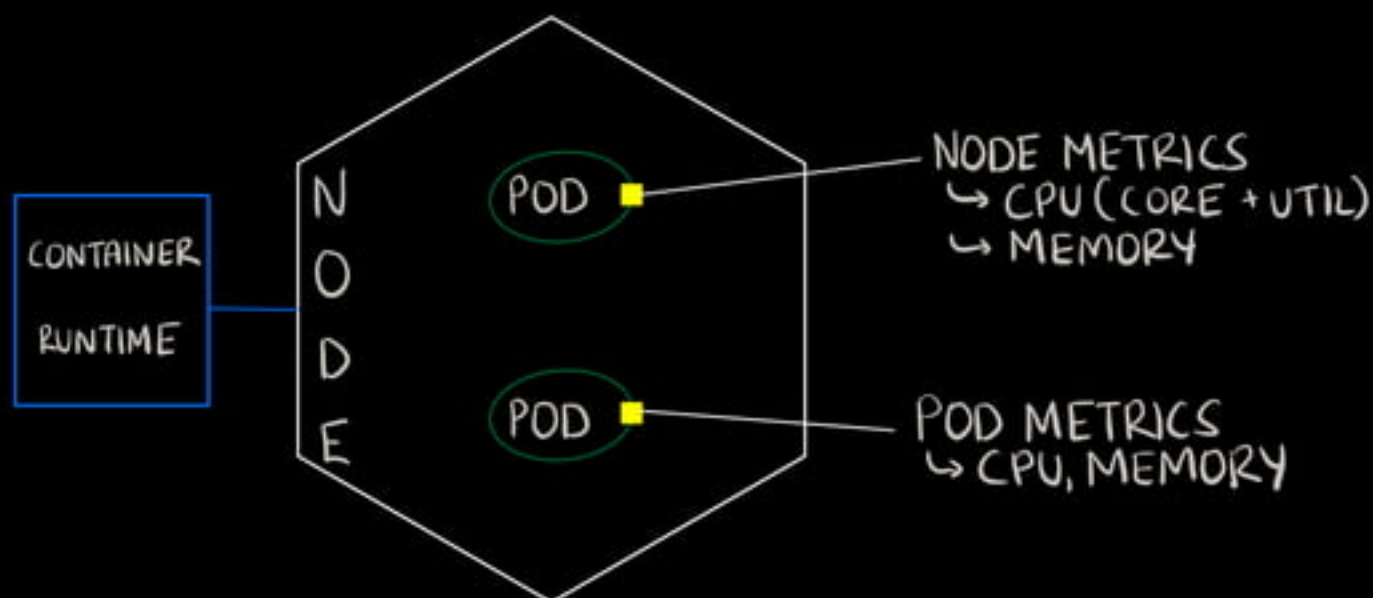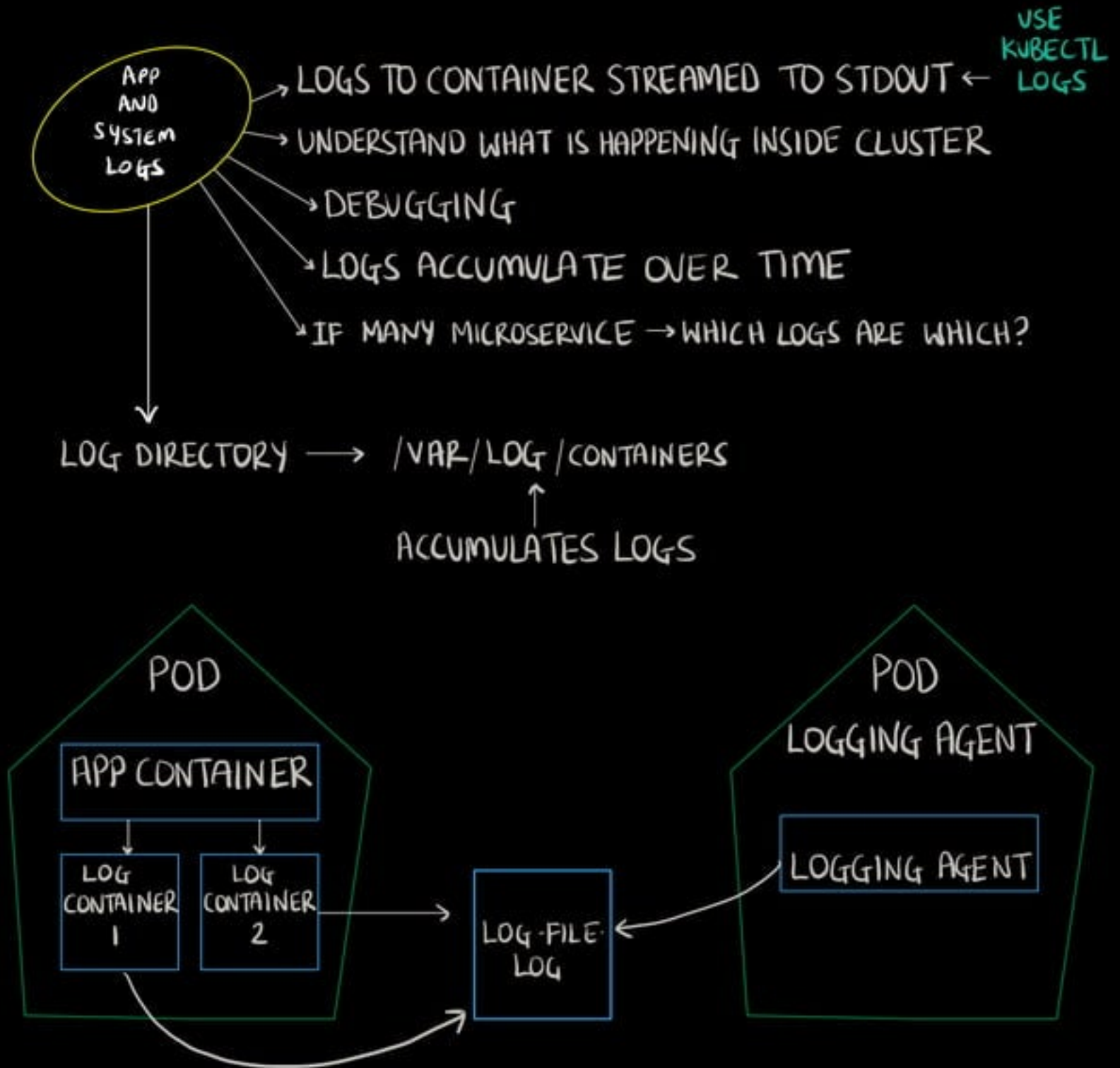
# Monitoring
# Cluster Components

# Monitoring the cluster components

The metric server allows you to collect CPU and memory data from the nodes and pods in your cluster



NODE METRICS
↳ CPU (CORE + UTIL)
↳ MEMORY

POD METRICS
↳ CPU, MEMORY

INSTALL METRIC SERVER →

KUBECTL TOP NODE ⟶ CPU|MEMORY FOR ALL THE NODES

KUBECTL TOP POD ⟶ CPU|MEMORY FOR ALL THE PODS

KUBECTL TOP POD --ALL-NAMESPACE ⟶ ALL NAMESPACE

KUBECTL TOP POD -N KUBE-SYSTEM ⟶ KUBE-SYSTEM NAMESPACE

KUBECTL TOP GROUP-CONTEXT --CONTAINERS → POD CONTAINERS

# Managing cluster component logs

**APP AND SYSTEM LOGS**
- → LOGS TO CONTAINER STREAMED TO STDOUT ← USE KUBECTL LOGS
- → UNDERSTAND WHAT IS HAPPENING INSIDE CLUSTER
- → DEBUGGING
- → LOGS ACCUMULATE OVER TIME
- → IF MANY MICROSERVICE → WHICH LOGS ARE WHICH?

LOG DIRECTORY ——→ /VAR/LOG/CONTAINERS

↑ ACCUMULATES LOGS

## POD

**APP CONTAINER**
- LOG CONTAINER 1
- LOG CONTAINER 2

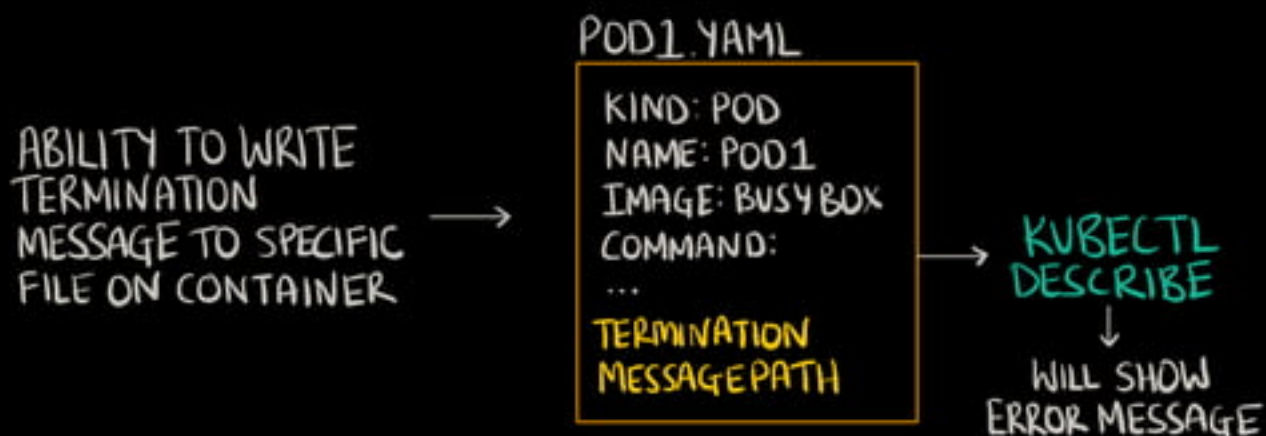→ LOG·FILE· LOG ←

## POD
### LOGGING AGENT

LOGGING AGENT

→ HAVE SIDECAR CONTAINER TO DO LOGGING SO YOU CAN ACCESS SPECIFIC LOGS

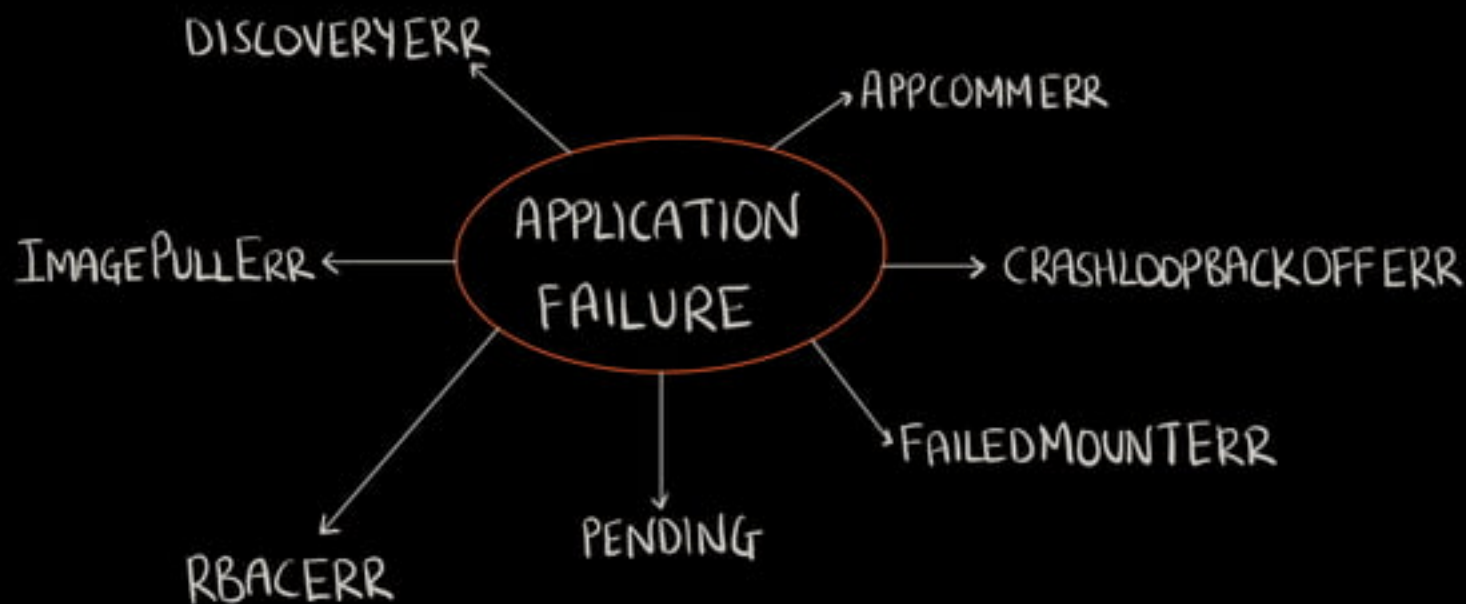→ ABLE TO ROTATE LOGS USING OTHER TOOLING → NO NATIVE

# Identifying Failures

# Troubleshooting Application Failure

ABILITY TO WRITE
TERMINATION
MESSAGE TO SPECIFIC
FILE ON CONTAINER

→

**POD1.YAML**

KIND: POD
NAME: POD1
IMAGE: BUSYBOX
COMMAND:
...
TERMINATION
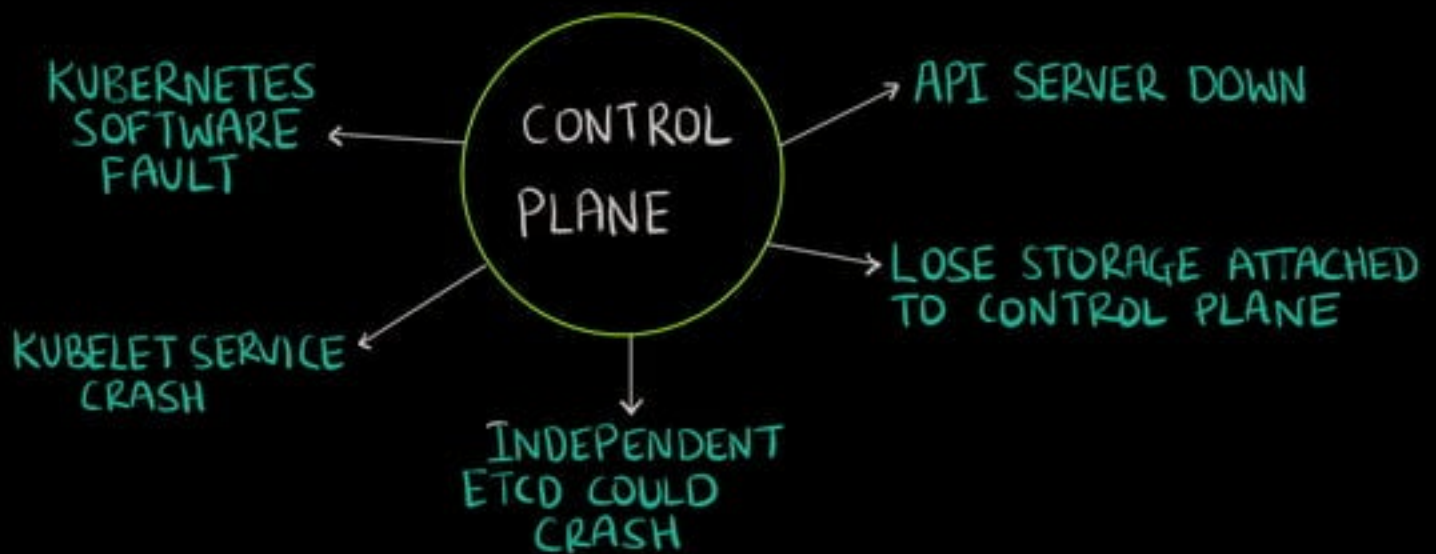MESSAGE PATH

→

KUBECTL
DESCRIBE
↓
WILL SHOW
ERROR MESSAGE

→ ONLY PARTICULAR FIELDS CAN BE CHANGED   I.E IMAGE

→ TO CHANGE OTHER FIELDS OF FAILED POD
↓
EXPORT CONFIGURATION
↓
MODIFY YAML I.E. CHANGE MEMORY REQUEST

DISCOVERYERR

APPCOMMERR

IMAGEPULLERR ←

**APPLICATION FAILURE**

→ CRASHLOOPBACKOFFERR

→ FAILEDMOUNTERR

RBACERR

PENDING

# Troubleshoot failures

**CONTROL PLANE**

- KUBERNETES SOFTWARE FAULT
- API SERVER DOWN
- LOSE STORAGE ATTACHED TO CONTROL PLANE
- KUBELET SERVICE CRASH
- INDEPENDENT ETCD COULD CRASH

→ VIEW THE EVENTS FROM CONTROL PLANE COMPONENTS

→ VIEW LOGS FOR CONTROL PLANE PODS

→ CHECK STATUS OF DOCKER SERVICE

→ CHECK STATUS OF KUBELET SERVICE

→ DISABLE SWAP

→ CHECK FIREWALLD SERVICE

→ VIEW KUBE CONFIG

**WORKER TROUBLESHOOTING**

- VIEW SYSLOG EVENTS
- VIEW KUBELET JOURNAL CTL LOGS
- GENERATE NEW TOKEN
- VIEW STATUS
- GET DETAILED INFO
- PING NODE
- IP ADDRESS OF NODE
- SSH INTO NODE

**NETWORK**

- ACCESS PODS DIRECTLY
- LOOKUP SERVICE VIA DNS
- CHECK IP TABLE RULES
- DNS RESOLVE CONF
- ENDPOINTS OF SERVICE
- CNI PLUGIN
- KUBERNETES SERVICE