



Kubernetes Workshop

Basic Concepts

Jan Harrie

Before we get started



- Technical material available on Github:

<https://github.com/NodyHub/k8s-ws>

```
git clone https://github.com/NodyHub/k8s-ws.git
```

- Make sure that you have an Minikube up and running
Else, use Vagrant and spin up a VM in Virtual Box with
https://github.com/NodyHub/k8s-ws/tree/master/00_prep/minikube
- Make sure that you have `kubectl` installed and configured

#whoami - Jan



- Freelancer – Security Consultant
- Former Security Analyst/Consultant/Pentester
- M.Sc. IT-Security TU Darmstadt
- Interests:
 - Orchestration Solution
 - Cloud
 - Gardening

#whoareyou?



- Expectation
- Background?
 - Dev
 - Ops
 - Networker
- Experience Docker/Kubernetes/PaaS

Re-Cap

Container



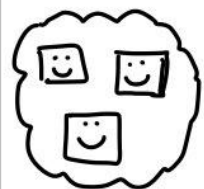
History

- Chroot circa **1982**
- FreeBSD Jails circa 2000
- Solaris Zones circa 2004
- Meosys - MetaClusters with Checkpoint/Restore 2004-05
- Linux OpenVZ circa 2005 (not in mainstream Linux)
- AIX WPARs circa 2007
- LXC circa 2008
- Systemd-nspawn circa 2010-2013
- Docker circa **2013**

JULIA EVANS
@b0rk

what's a container?

a Linux container
is a group of processes



Container

We have our own
filesystem but
we're still just
regular processes!

Linux containers are
isolated from other processes

they can have their own:

- users
- network namespace
- filesystem
- process IDs
- memory / CPU limits

Kernel features that
isolate Linux containers

cgroups

namespaces

capabilities

seccomp-bpf

there are many ways
to run Linux containers

runc

systemd-nspawn

LXC

Docker

Docker
uses runc
under the
hood

your own homegrown
bash script

and containers can be
set up in different ways



container 1

I have my own
filesystem!

I don't! I
have my system
calls restricted!



container 2

extra confusion:

"container" sometimes means
"lightweight VM"

Fargate and kata Containers
are actually VMs and not
Linux containers (they don't
share a kernel with other
containers)

Build a container:

00_prep/k8s-ws/

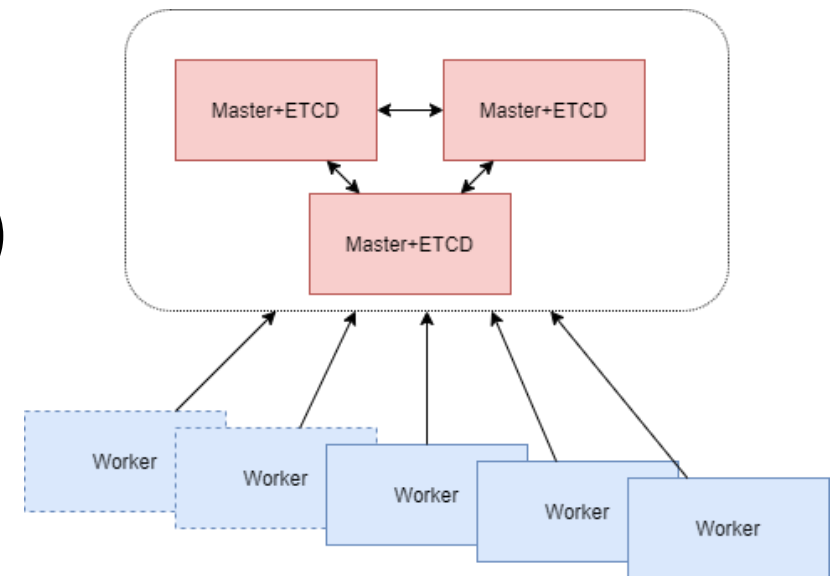
Exercise

Kubernetes

Architecture

What is Kubernetes?

- Widespread Container Orchestration solution
- Similar concepts to Docker Swarm
 - Uses RAFT for master node
 - Secrets
 - Services / Deployments
- Former Borg¹ (originally developed by Google)
- Uses etcd KV-Store



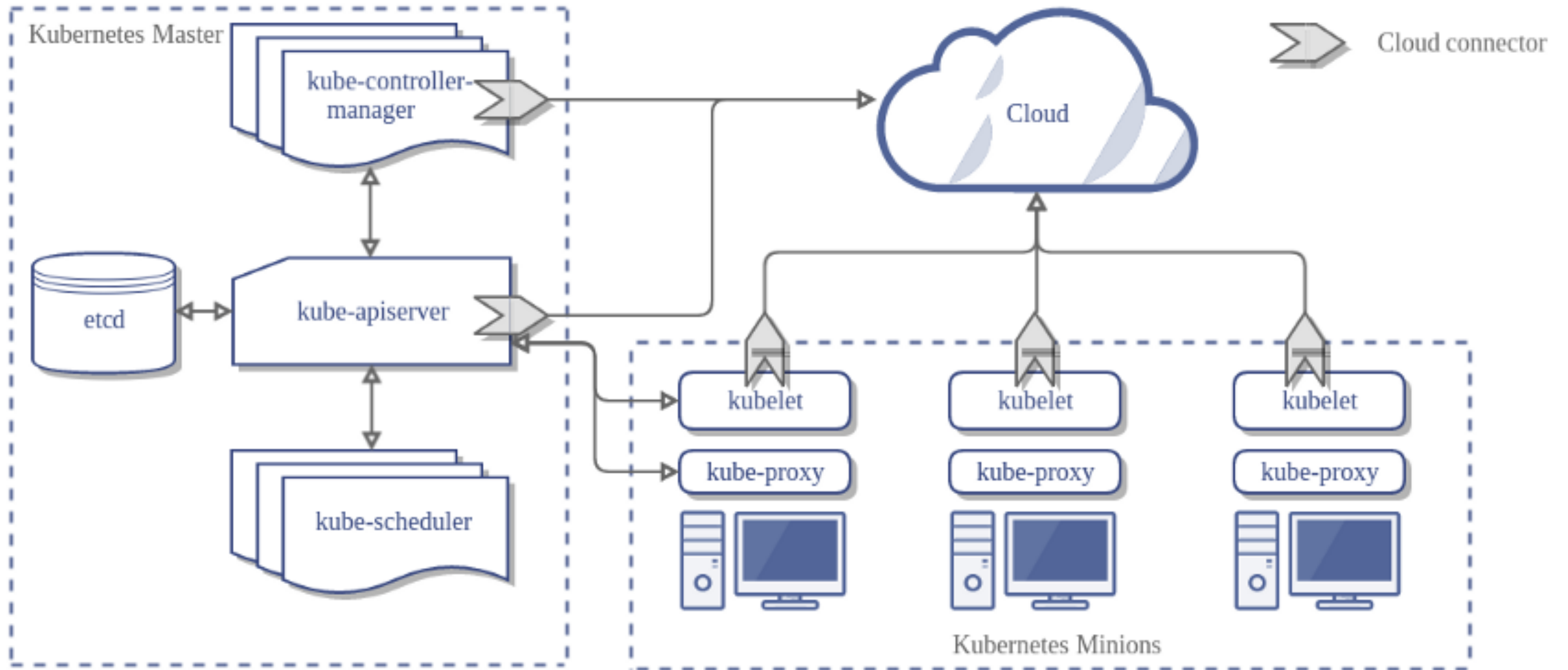
¹: <https://research.google.com/pubs/pub43438.html>

What is Kubernetes?

- Manager nodes handle cluster management tasks
 - Cluster state maintenance
 - Scheduling of services
 - API endpoints
- Worker nodes provide executes the wotkload

Kubernetes Clusters

- Builds on the Raft Consensus Algorithm
- Manager nodes keep the cluster state consistent in etcd
- In case of failure:
 - majority of nodes needs to agree on values
 - $(N-1)/2$ failures tolerated, otherwise no more requests are processed



Kubernetes CLI (kubectl)

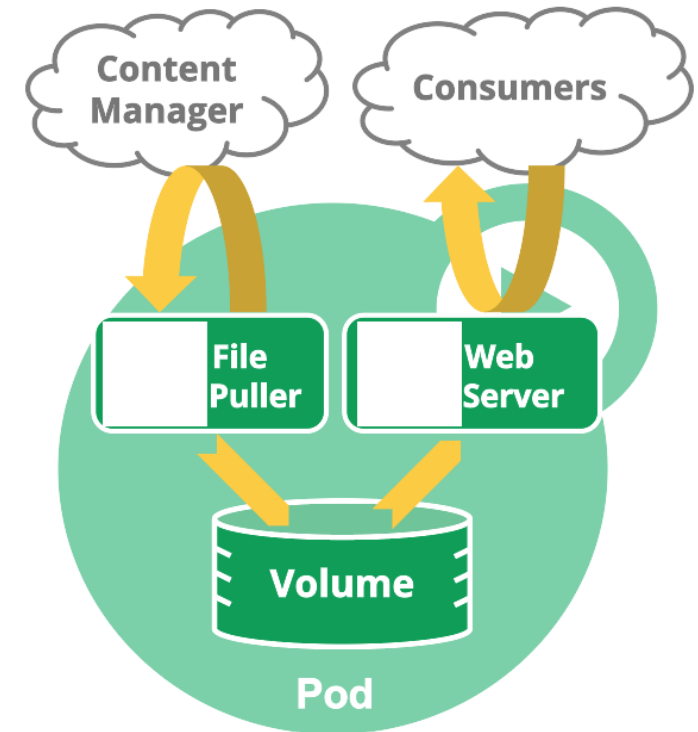
- kubectl get
 - Get a resource
 - i.e. kubectl get nodes
- kubectl describe
 - Describe a resource
- kubectl apply
 - Apply a resource file
- kubectl create / kubectl run / kubectl expose
 - Create / manage services / deployments

Kubernetes

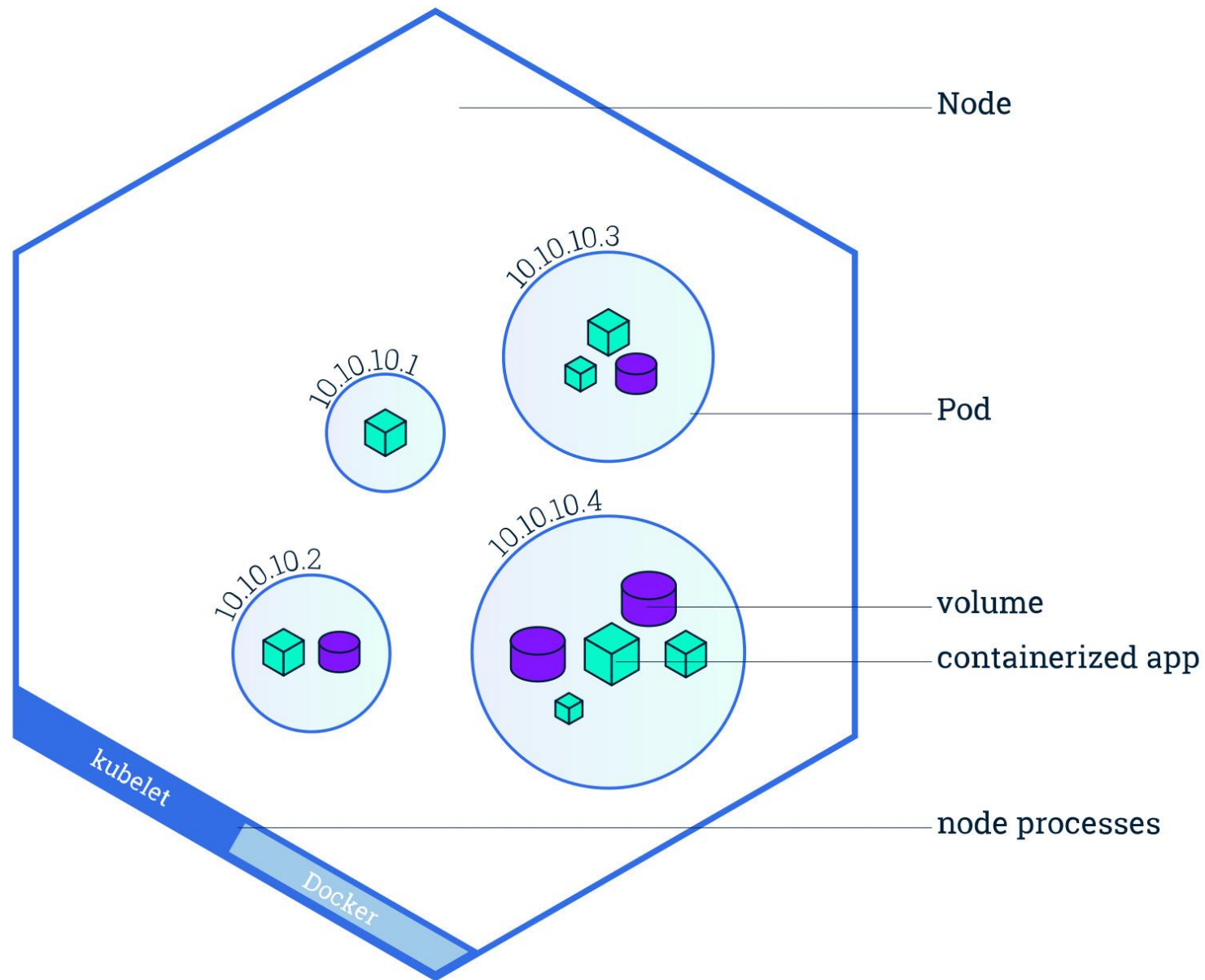
Pod

Pod

- Smallest unit in the Kubernetes ecosystem
- A pod is a group (≥ 1) of containers
- Are co-scheduled and share resources
- Shared IP/Ports
- C-to-C communication: IPC
- Best practice: 1 Container per Pod
- Health checks for containers



Pods



Start a pod

01_getting-started/00-pod/

Exercise

Pod Lifecycle (1/2)

- Status:
 - Pending: Created, but not scheduled yet
 - Running: Scheduled
 - Succeeded: All containers terminated successfully
 - Failed: At least one container terminated with an error
 - Unknown: API server was not able to query state
 - CrashLoopBackoff: Container failed to start, trying it again
 - ImagePullBackOff: Image cannot be pulled

Pod Lifecycle (2/2)

- Container probes
 - livenessProbe
 - readinessProbe
 - Implemented by
 - ExecAction
 - TCPSocketAction
 - HTTPGetAction
- On stop
 - PreStop hook called
 - Containers are killed

Kubernetes

Deployment

Workloads

- Pod
- ReplicaSet (former ReplicationController)
- Deployment
- StatefulSet
- DaemonSet
- (Cron-)Job (Batch-style workload)

ReplicaSet

- Successor of ReplicationController
- Ensures that no. of pod replicas are running
- Mostly used in combination with Deployments

Deployment

- ReplicaSet with control
- Regularly check status of pods
- Update & rollback possible
 - Rolling updates
 - Blue-Green deployments

DaemonSet

- Bypasses the default scheduler
- Run a pod on every node
- Used to i.e. bootstrap
 - Network plugins make use of it

Job

- Creates one or more pods
- Make sure a specified number of them terminates successfully => Job done
- i.e. simple job which has to run once
 - Job will start new pod, if first pod fails or is deleted (i.e. due to node hardware failure)
- Multiple pods possible
- Cronjob: timed Job

Create a Deployment

01-deployment

Exercise

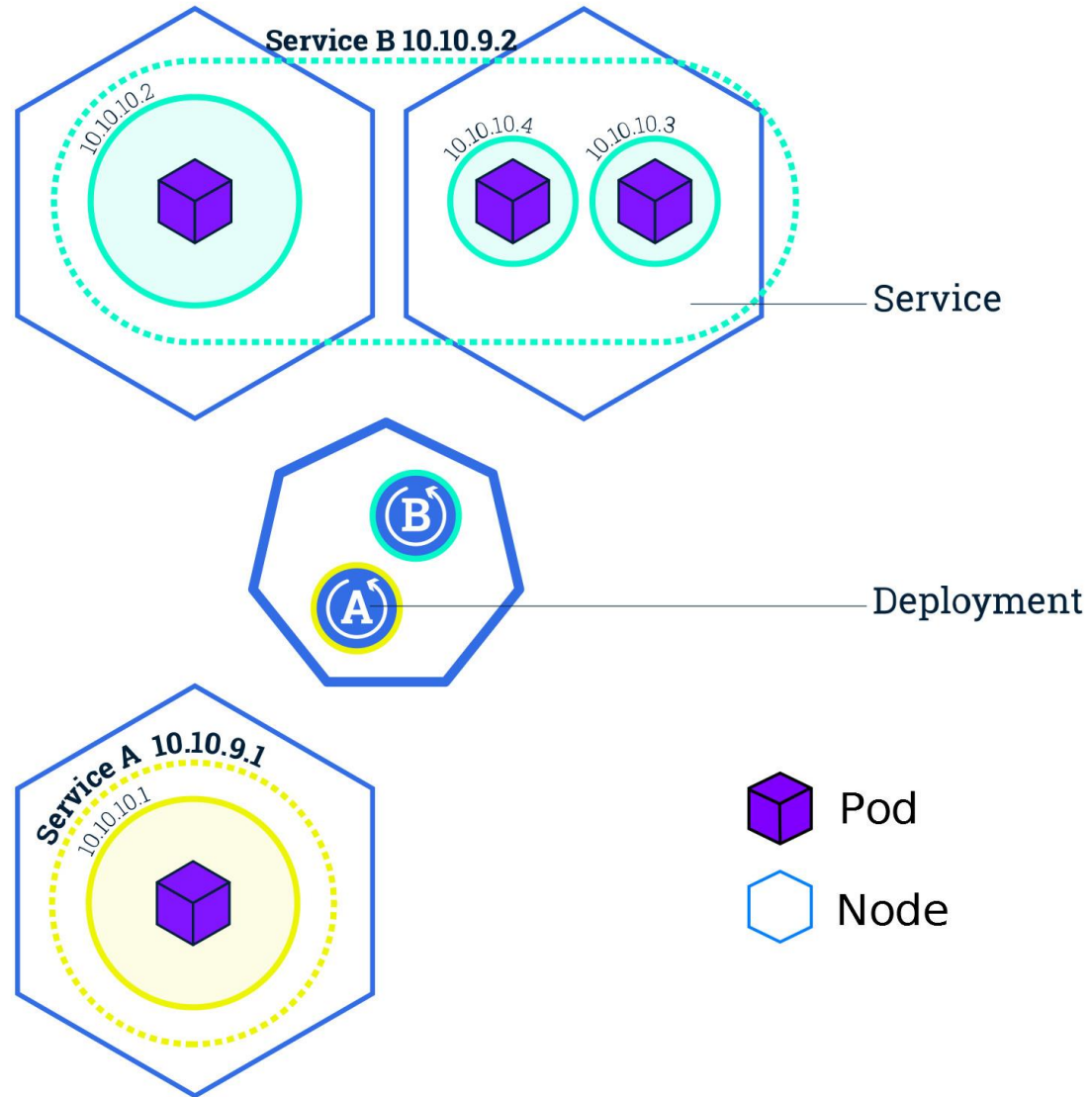
Kubernetes

Service

Services

- Problems
 - Pods are mortal
 - Multiple Pods for one service exist
 - No access from outside
- A Service defines a set of Pods and a policy (i.e. ports)
 - Including a VIP
- Service types
 - ClusterIP (only internal, default)
 - NodePort (binds port on each node)
 - LoadBalancer (use a cloud provided load balancer)
 - ExternalName (CNAME entry)

Services



Service Discovery (1/3)

- Service Discovery is build up on DNS
- CoreDNS is used as backend
- Every service receives a unique DNS name
 - i.e. service foo in namespace bar
 - Pod in same namespace: foo
 - Pod in other namespace: foo.bar
 - Can also be enabled for pods
 - 1-2-3-4.default.pod.cluster.local

Service Discovery (2/3)

- A Record
 - <service>.<ns>.svc.<zone>.
 - kubernetes.default.svc.cluster.local.
- SRV Record (service record)
 - _<port>._<proto>.<service>.<ns>.svc.<zone>.
 - _https._tcp.kubernetes.default.svc.cluster.local.
- PTR Record (reverse DNS)
 - <a>..<c>.<d>.in-addr.arpa.
 - 1.0.3.10.in-addr.arpa.

Service Discovery (3/3)

```
vagrant@ubuntu-bionic: ~/git/k8s-basic-ws/01_getting-started/00...  
File Edit View Search Terminal Help  
ubuntu-bionic [../01_getting-started/00-pod]% kubectl exec -it k8s-ws -- sh  
/app # cat /etc/resolv.conf  
nameserver 10.96.0.10  
search default.svc.cluster.local svc.cluster.local cluster.local fritz.box  
options ndots:5  
/app #
```

Create a Service

02-service

Exercise

Kubernetes

Secrets & Configmaps

ConfigMaps

- Pass information to workloads (configuration)
- Decoupling of configuration from images
- Consume by ENV or mounted volumes

Secrets

- Similar to ConfigMaps
- Pass information to workloads (secrets)
- Decoupling of secrets from images
- Consume by ENV or mounted volumes

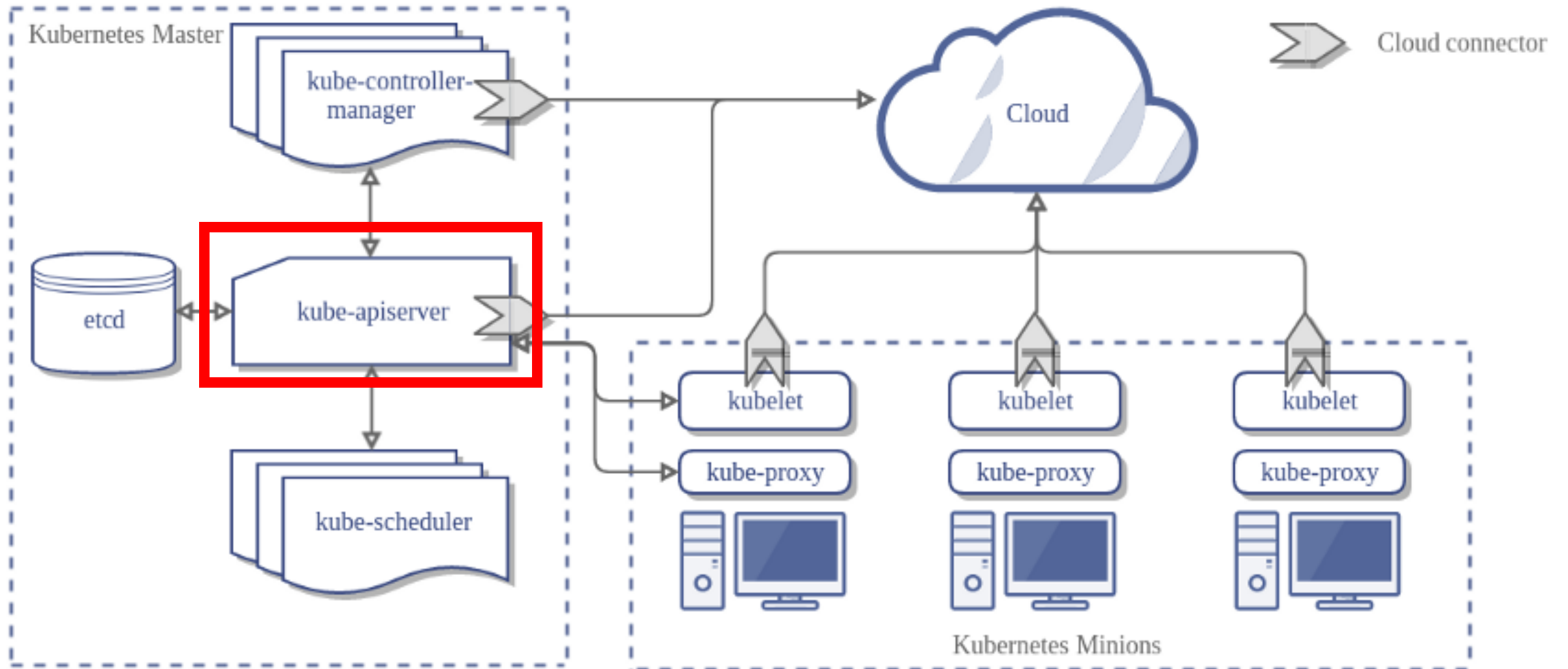
Secrets & Configmaps

03-secrets-configmaps

Exercise

Kubernetes

Master Components, etcd & Yaml Files

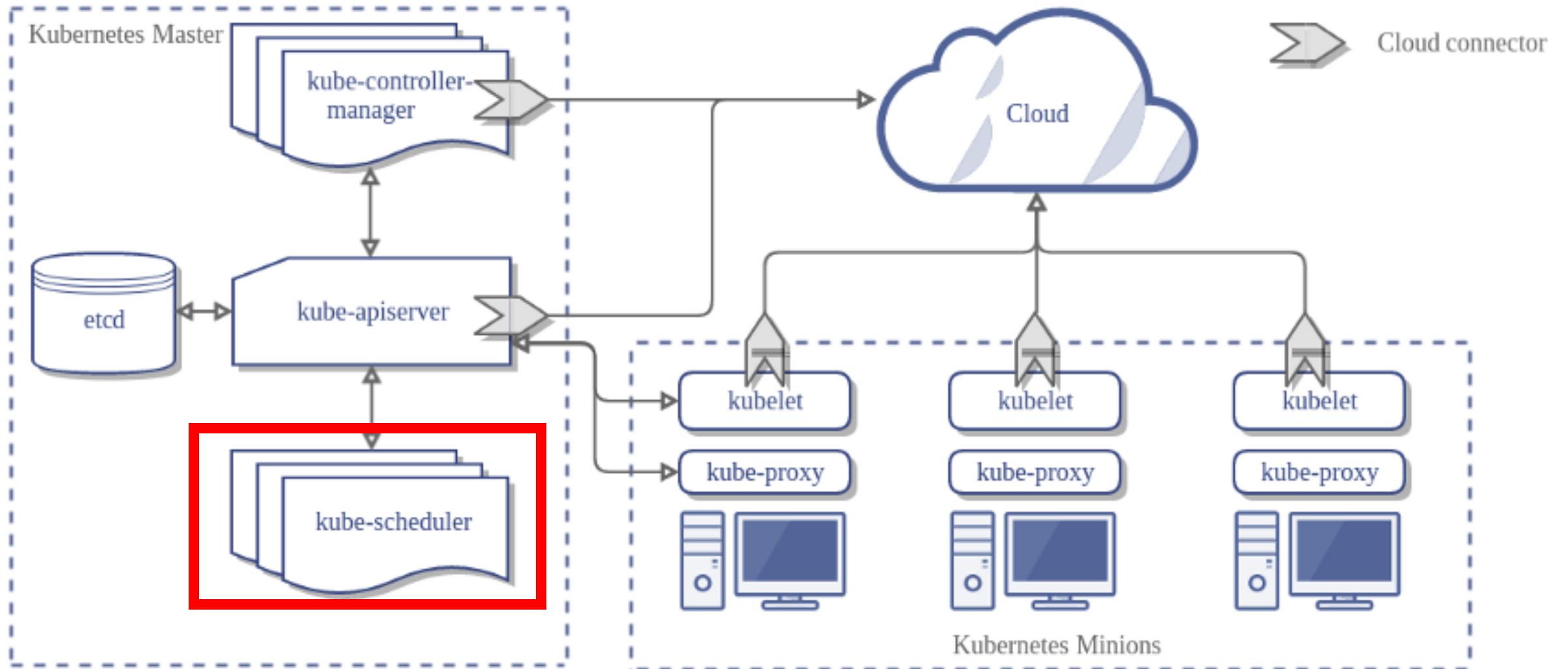


Kubernetes API

- REST API / CRUD
- HTTP API with JSON
- Versioned
 - /api/v1
 - /apis/extensions/v1beta1
- Base path
 - /apis/batch/v1/namespaces/....
 - /apis/extensions/...
- Core API (historically not /apis/core/v1)
 - /api/v1
- Custom Resource Definitions (CRD) possible

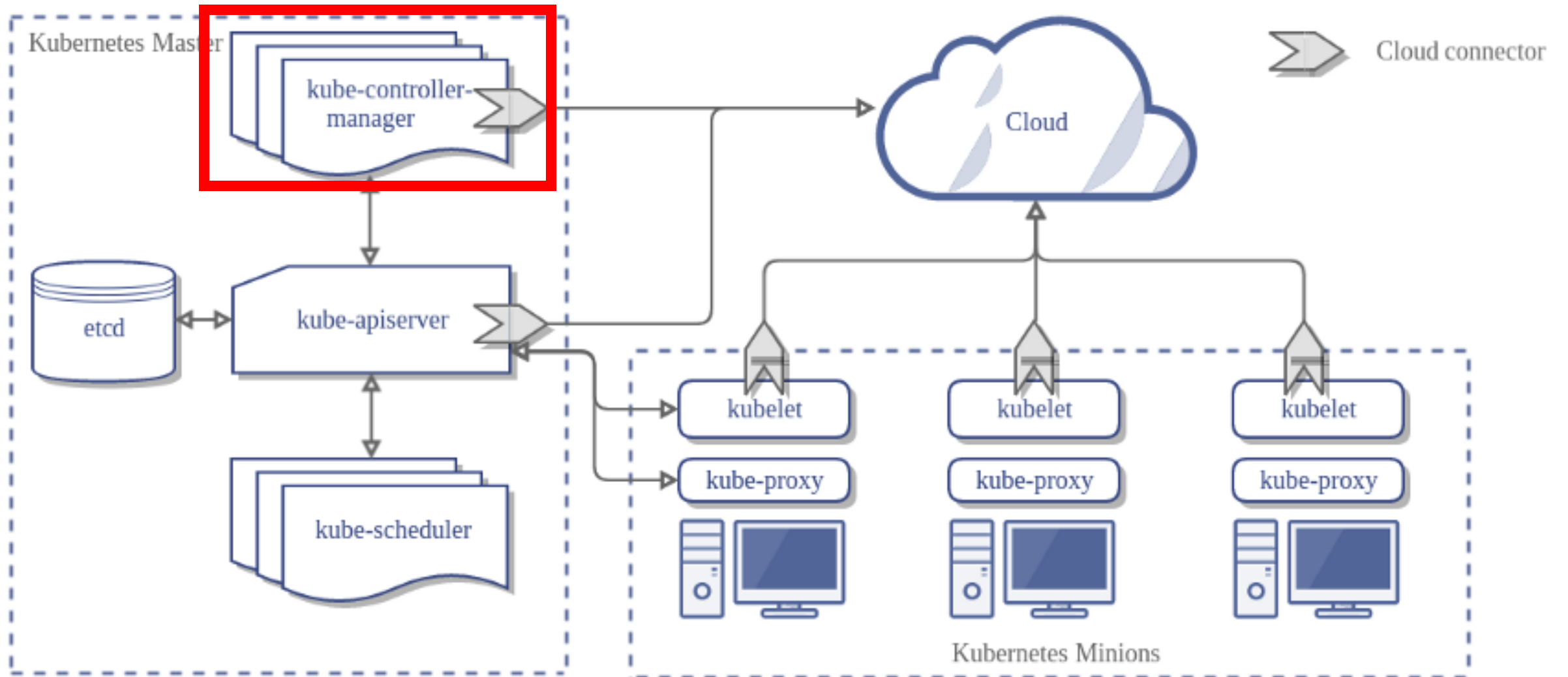
Kubernetes API (Objects)

- apiVersion
- kind (Object type)
 - Persistent entity (Deployment, Service, Pod, ...)
 - List (collection of resources)
 - Special purpose (i.e. bindings, status)
- metadata (labels, name)
- spec (desired state of the object)
 - “Kind”-specific data



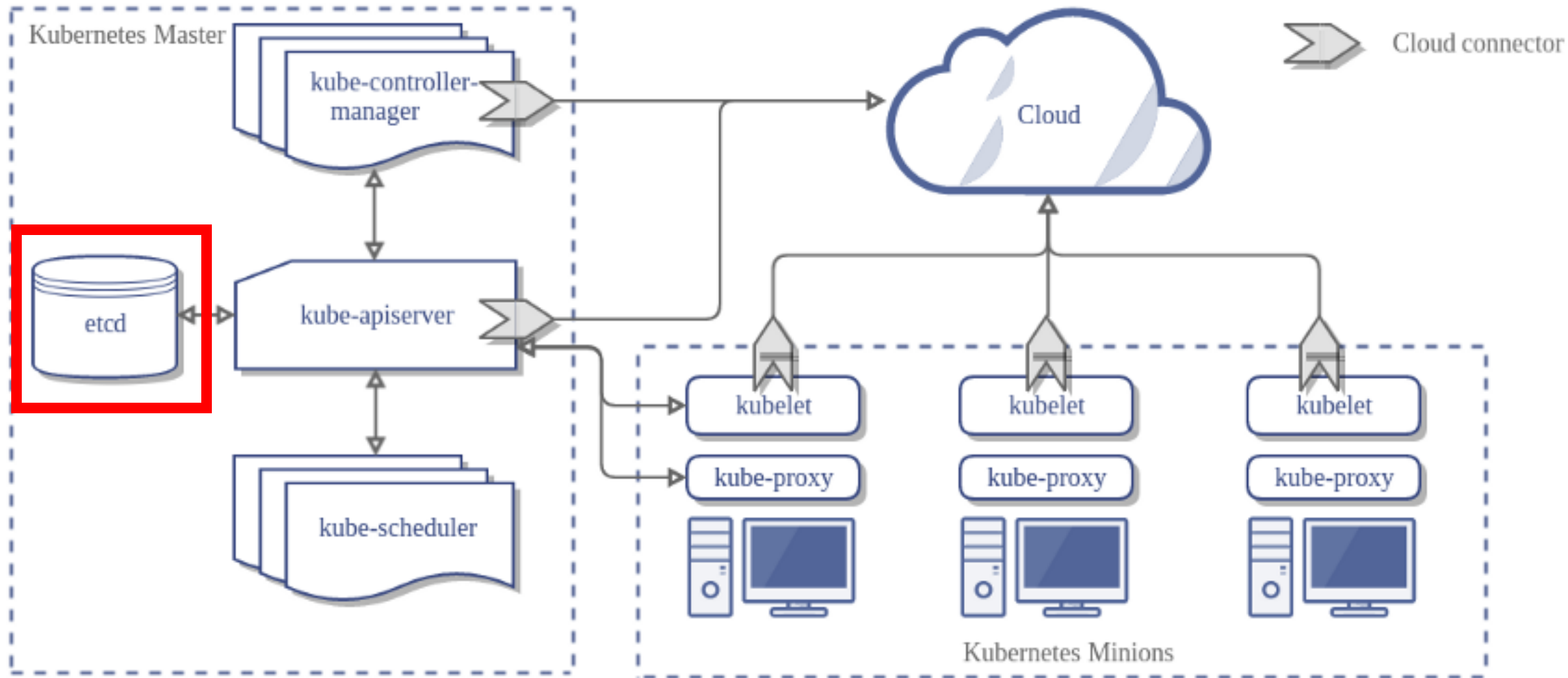
kube-scheduler

- Binds unscheduled pods to nodes
- Pluggable: multiple (custom) schedulers possible



Controller Manager

- Single binary on master which contains controllers
 - Endpoint Controller (manages Endpoints)
 - Node Controller (manages nodes)
 - Replication Controller (maintain no. of pods)
 - Service Account and Token Controller
 - Cloud Controller (cloud integration)



etcd

- Distributed persistence layer (key-value storage)
- Stores the cluster state
- Distribution / Replication
 - Leader elections
 - Fault tolerance
 - Distributed locks
- v2 vs v3
- Authn / Authz (disabled by default)
- *etcdctl* to control etcd



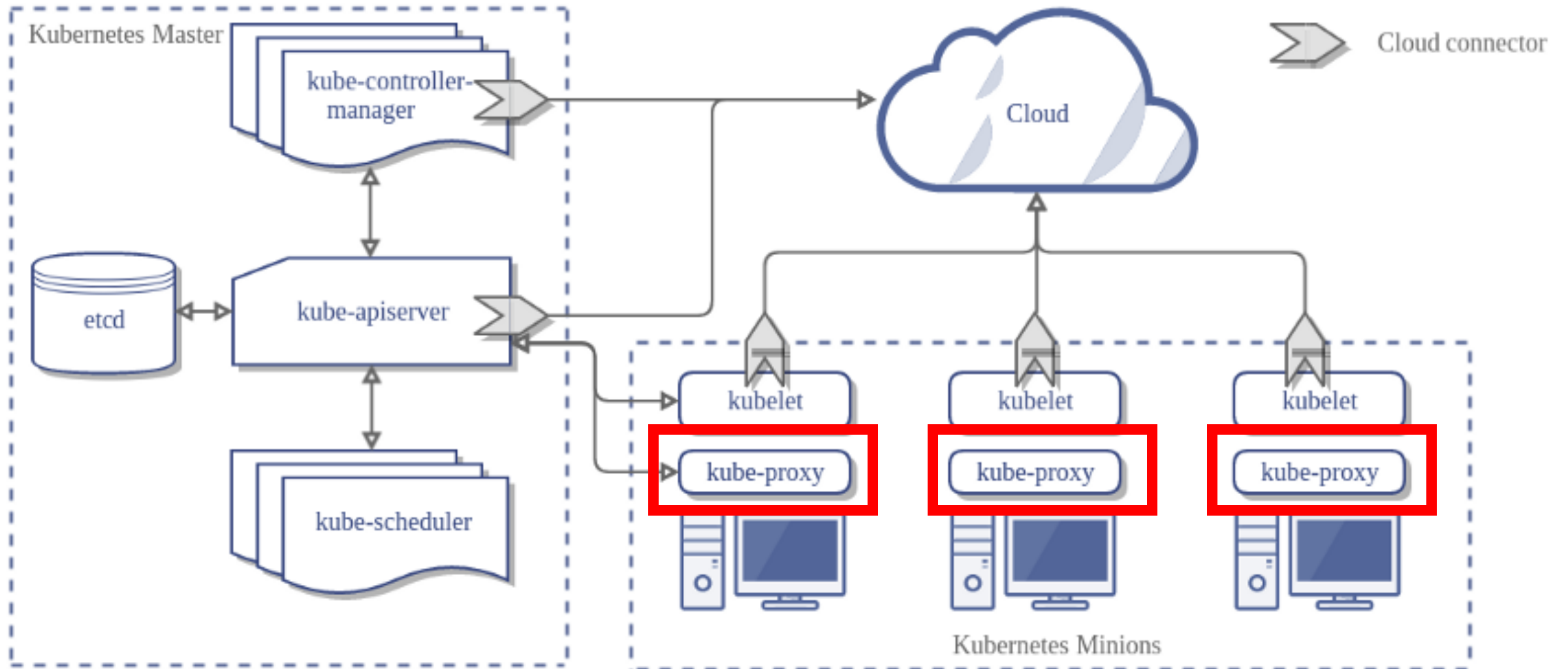
Yaml Files

04-yaml

Exercise

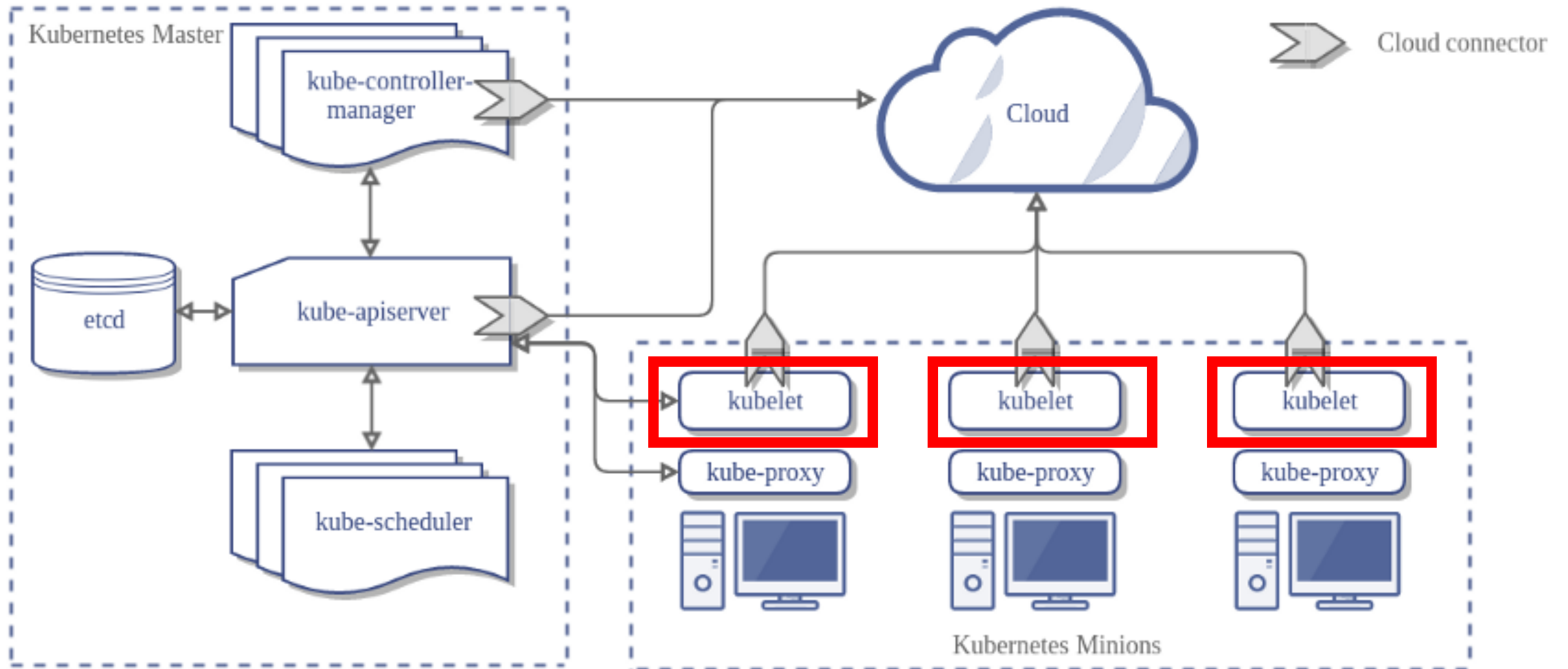
Kubernetes

Node Components



kube-proxy

- Route packages to pods
- Load balancer for pods
- Maintain network rules (iptables)
- “Services” for the nodes



kubelet

- Manages pods, containers, volumes, secrets
- Agent which runs on every computing node (nodes + sometimes master)
- Authentication for kubelet (port 10250)
 - X509 client certificate
 - WebHook
 - Calls *TokenReview* API

Manage the environment

05-env

Exercise

Kubernetes

Labels & Selectors

Labels and Selectors

- Labels are key / value pairs that are attached to objects (i.e. Pods)
 - i.e. environment (dev, testing, prod), release (stable, dev), tier (frontend, backend, cache)
- Selectors are used to query objects based on labels
 - Restart all pods with label production
 - Required for Service, Ingress, LoadBalancer, NetworkPolicies

Manage the environment

06-link-configmap

Exercise

Manage the environment

07-link-secret

Exercise

Kubernetes

Ingress

Ingress

- L4-L7 decisions
 - Defines rules
 - Redirects traffic to Services
- “Service load balancer”
- SSL termination support
- Managed by Ingress Controller
 - Part of *kube-controller-manager* binary
 - Supports GCE and nginx controllers

Manage the environment

08-ingress

Exercise

Kubernetes

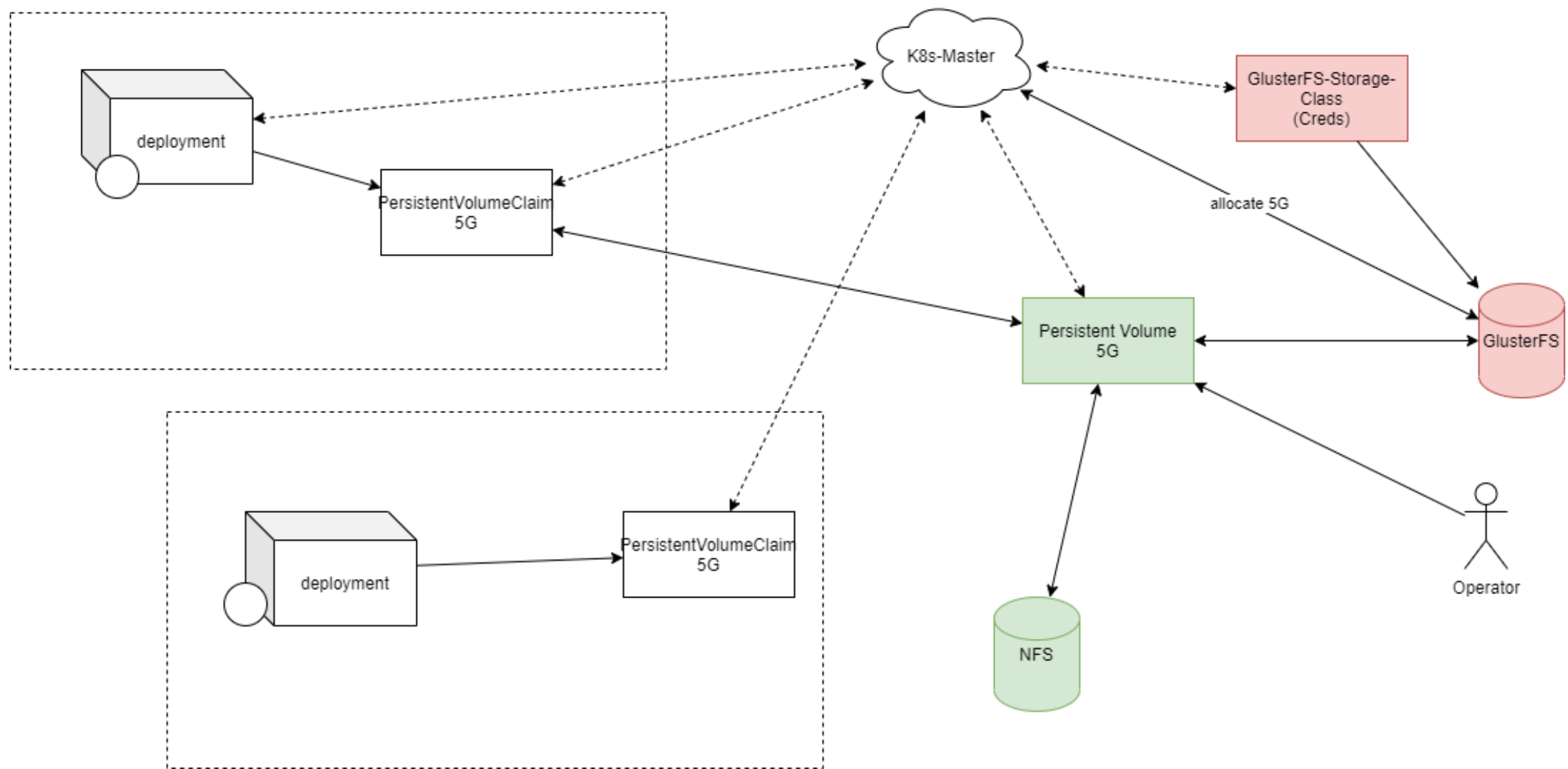
Persistence

Volumes

- Per default data is ephemeral
 - Data lifetime == Pod lifetime
 - i.e. used for caching
- Persistence
 - Data lifetime != Pod lifetime
 - Concept of Volumes
 - Core: just a directory with specified backend
 - Types: emptyDir, hostPath, nfs, iscsi, glusterfs

PersistentVolumes API

- Abstraction of details how storage is provided and consumed
- PersistentVolume (PV)
 - Storage provisioned by administrator
 - Labeled
- PersistentVolumeClaim (PVC)
 - Request for storage by user
 - Request specific class, size and access mode
- StorageClasses (i.e. SSD, HDD)
 - provisioner (i.e. AWS Block storage, iSCSI, NFS)
 - parameter (i.e. type, zone, ...)
 - reclaimPolicy (i.e. Delete or Retain)



Manage the environment

09-pvc

Exercise

Thanks for
your Attention

- Any Further Question?

Jan Harrie

Mail: jharrie@harrie-consulting.de

Twitter: [@NodyTweet](https://twitter.com/NodyTweet)

Blog: <https://blog.nody.cc/>

GitHub: <https://github.com/NodyHub>

Dockerhub: <https://hub.docker.com/u/nodyd>