

```
<!-- TOC -->
```

```
* [General](#general)
* [Create pods](#create-pods)
* [Test a pod](#test-a-pod)
* [Pods, containers and storage](#pods-containers-and-storage)
* [Create other resources](#create-other-resources)
* [Update resources](#update-resources)
* [Debugging](#debugging)
* [Delete / replace resources](#delete--replace-resources)
* [Secrets for ServiceAccount](#secrets-for-serviceaccount)
* [Networking, services, DNS](#networking-services-dns)
* [Cluster](#cluster)
* [Helm](#helm)
* [etcd](#etcd)
```

```
<!-- TOC -->
```

General

- In Mousepad, prepare these setup commands:

```
cp $HOME/.bashrc $HOME/.bashrc_backup # create a backup file, just in case...
git clone https://github.com/ghahitette/kubernetes
cd kubernetes/
chmod +x bashrc_append.sh
./bashrc_append.sh
```

- To assume root privileges: ``sudo -i`` (only applies to ACG server setup?)
- Use ``grep -A2 Mounts`` to show two lines after the line matching ``Mounts``
- Select the ``acgk8s`` cluster to interact: ``k config use-context acgk8s``
- Watch pods / deployments / jobs: ``k get pods -w`` / ``k get deployments -w`` / ``k get jobs -w``
- Repeat command every n seconds, example: ``watch -n 2 kubectl get pods``
- Check all resources [in all namespaces]: ``k get all [-A]``
- List k8s "internal" pods, sorted by node name: ``k get pods -n kube-system --sort-by .spec.nodeName``
- List of all / all namespaced resources: ``k api-resources [--namespaced] [-o (wide | name)]``
- API e.g. for pod manifests : ``k explain pods[child1.child2] | more``

k8s architecture

- Control Plane: components managing the cluster itself globally, usually run on dedicated controller machines.
 - ``kube-api-server``: the primary interface to the control plane and the cluster itself.
 - ``kube-scheduler`` selects available nodes on which to run containers.
 - ``kube-controller-manager`` runs a collection of multiple controller utilities in a single process
 - ``etcd``: the HA backend data store for the Kubernetes cluster.
 - ``cloud-controller-manager``: interface between Kubernetes and various cloud platforms (optional).
- Nodes: the machines where the containers run. A cluster can have any number of nodes
 - ``kubelet``: the Kubernetes agent that runs on each node.
 - Ensures that containers are run on its node, as instructed by the control plane.
 - Reports container data (e.g. status) back to the control plane.
 - ``container runtime``: runs containers (e.g. Docker, containerd)!!! not built into Kubernetes
 - ``kube-proxy`` is a network proxy, provides networking between containers and services in the cluster.

Cluster

- Get Services | pods IPs range: ``k cluster-info dump | grep -m 1 (service-cluster-ip-range | cluster-cidr)`` or describe the ``kube-controller-manager`` pod
- Static Pod = a Pod managed directly by ``kubelet`` on a node, not by the K8s API server; can run even without a K8s API server present; created from YAML manifest files in ``/etc/kubernetes/manifest/`` (by default)
- ``kubelet`` creates a mirror Pod for each static Pod, sharing the status of the static Pod to the k8s API
- Networking (CNI plugin) is configured on control plane node(s) under ``/etc/cni``, e.g. ``/etc/cni/net.d``
- To temporarily stop ``kube-scheduler``, log onto the control plane node, move its YAML manifest file (e.g. to ..) and restart ``kubelet``
- To manually schedule a Pod on a node, set ``pod.spec.nodeName``, and not ``pod.spec.nodeSelector`` (works even if ``kube-scheduler`` is not running)
- Pod termination: when available cpu or memory resources on the nodes reach their limit, first candidates for termination are Pods using more resources than they requested (by default containers without resource requests/limits set).
- A DaemonSet ensures that all (or some) Nodes run a copy of a Pod (e.g. for network plugins, cluster storage, logs collection, node monitoring)
 - To create a DaemonSet, create a Deployment YAML file with ``kubectl`` and modify (remove ``replicas``, ``strategy`` and add...)
 - As nodes are added to / removed from the cluster, Pods are added / garbage collected. Deleting a DaemonSet will clean up the Pods it created.
 - ``spec.selector`` is a pod selector, immutable and must match ``spec.template.metadata.labels``
 - By default, DaemonSet pods are created and scheduled by the DaemonSet controller, not the Kubernetes scheduler. ``ScheduleDaemonSetPods`` allows you to schedule DaemonSets using the default scheduler instead of the DaemonSet controller, by adding the ``NodeAffinity`` term to the DaemonSet pods, instead of the ``spec.nodeName`` term.
- Drain a node: ``k drain [--ignore-daemonsets --force] <node name>``
 - The ``kubectl drain`` subcommand on its own does not actually drain a node of its DaemonSet pods: the DaemonSet controller (part of the control plane) immediately replaces missing Pods with new equivalent Pods.
 - The DaemonSet controller also creates Pods that ignore unschedulable taints, which allows the new Pods to launch onto a node that you are draining.
- Resume scheduling ``new pods`` onto the node: ``k uncordon <node name>``
- In a cluster built with ``kubeadm``:
 - Check the status of cluster components (e.g. kube-apiserver): check the status of (static) Pods in the ``kube-system`` Namespace (kube-apiserver is not set up as a systemctl service).
 - Find logs for the Kubernetes API Server: ``k logs -n kube-system <api-server-pod-name>`` (the ``/var/log/``

kube-apiserver.log` log file is not available on the host since the API Server runs in a static Pod).

- Find kubelet logs: `sudo journalctl -fu kubelet` (kubelet runs as a standard service).
- Investigate DNS issues: check the DNS Pods in the `kube-system` Namespace.
- Build / upgrade `kubeadm` clusters: [link](CKA%20training/Build%20kubeadm%20clusters.md) / [link](CKA%20training/Upgrade%20kubeadm%20clusters.md)
- Certificates: for Kube API server certificates, ssh to control plane node and `kubeadm certs (check-expiration | renew)`
- Certificates: for kubelet client/server certificates, ssh to the node, check `--cert-dir` parameter for the kubelet or `/etc/systemd/system/kubelet.service.d/10-kubeadm.conf` and `openssl x509 -noout -text -in /var/lib/kubelet/pki/kubelet-client-current.pem | grep Issuer` or `openssl x509 -noout -text -in /var/lib/kubelet/pki/kubelet.crt | grep "Extended Key Usage" -A1`
- Scenario: broken `kubelet` on a node (showing as `NotReady`) with `/usr/bin/local/kubelet` not found error: ssh to the node, find the kubelet service config file with `systemctl status kubelet`, find the kubelet location with `whereis kubelet`, modify config file `/etc/systemd/system/kubelet.service.d/10-kubeadm.conf` to fix path to `/usr/bin/kubelet` and `systemctl daemon-reload && systemctl restart kubelet`
- Scenario: un-initialised (`/etc/kubernetes/kubelet.conf: no such file or directory`), outdated node needing to join the cluster: ssh on node, `apt install kubect1=1.26.0-00 kubelet=1.26.0-00` then ssh on control plane, `kubeadm token create --print-join-command` then `sudo kubeadm join ...` command from node
- Scenario: curl the k8s API from a test pod using a ServiceAccount `secret-reader`: retrieve the ServiceAccount token `TOKEN=\$(cat /var/run/secrets/kubernetes.io/serviceaccount/token)` and `curl -k https://kubernetes.default/api/v1/secrets -H "Authorization: Bearer \${TOKEN}"`
- To use encrypted https connection: `CACERT=/var/run/secrets/kubernetes.io/serviceaccount/ca.crt` and `curl --cacert \${CACERT} https://kubernetes.default/api/v1/secrets -H "Authorization: Bearer \${TOKEN}"`

Create pods

- Create an nginx pod: `k run my-pod --image=nginx [--port=80] [--labels app=b]`
- Create a busybox pod: `k run my-pod --image=busybox \$do --command -- sh -c "sleep 1d"`
 - Command YAML syntax example: `command: ['sh', '-c', 'sleep 1; done']`
- Create a throw-away, interactive pod with busybox | netshoot: `k run my-pod --image=(busybox | nicolaka/netshoot) --restart=Never --rm -ti`

Test resources

- Test a pod with a command: `k exec my-pod [-c my-container] (-- env ... | -- cat ...)`
- Test a pod interactively: `k exec my-pod [-c my-container] -ti -- sh`
- Test RBAC: `k auth can-i`, example: `k auth can-i create configmap --as system:serviceaccount:project-hamster:processor`

Pods, containers and storage

- Startup probes: run at container startup and stop running once they succeed; very similar to liveness probes (which run constantly on a schedule); useful for legacy applications that can have long startup times.
- Readiness probes: used to prevent user traffic from being sent to pods that are still in the process of starting up (e.g. pod STATUS = Running but READY = "0/1")
 - Example: for a service backed by multiple container endpoints, user traffic will not be sent to a particular pod until its containers have all passed readiness checks.
- Pod's restart policy: Always (by default), OnFailure (restarted only if error code returned), and Never.
- Pod with InitContainer(s) will show "Init(0/n)" in their STATUS during initialisation
- Set the node name where a Pod is running as environment variable:

```
spec:
  containers:
    - image: nginx:1.17.6-alpine
      env:
        - name: MY_NODE_NAME
          valueFrom:
            fieldRef:
              fieldPath: spec.nodeName
  ...
```

- Volumes:
 - A PersistentVolume's `persistentVolumeReclaimPolicy` determines how the storage resources can be reused when the PersistentVolume's associated PersistentVolumeClaims are deleted:
 - `Retain`: Keeps all data. An administrator must manually clean up the data to prepare the storage resource for reuse.
 - `Delete`: Deletes the underlying storage resource automatically (cloud only).
 - `Recycle`: Automatically deletes all data, allowing the PersistentVolume to be reused.
 - `allowVolumeExpansion` property of a **StorageClass**, if set to false (per default), prevents from resizing a PersistentVolumeClaim.
- Kill the `containerd` container of the `kube-proxy` Pod on a given node: ssh to the node and

```
crictl ps | grep kube-proxy
crictl stop 1e020b43c4423
crictl rm 1e020b43c4423 # kubelet will restart the container with a new ID
...
```

- To read container logs from worker node using `crictl` and write them to exam server: `ssh cluster1-node2 'crictl logs b01edbe6f89ed' &> /opt/course/17/pod-container.log` # The &> in above's command redirects both the standard output and standard error`

Create other resources

- Create a job with `k create job my-job --image=busybox \$do > job.yml -- sh -c "sleep 2 && echo done"` then check pod execution (no such thing as starting a Job or CronJob!)
- Create a ConfigMap from a file, with a specific key: `k create configmap my-cm --from-file=index.html=/opt/course/15/web-moon.html`
- Create a secret (with implicit base64 encoding): `k create secret generic my-secret --from-literal user=test --from-literal pass=pwd`
- Create an nginx deployment: `k create deployment my-dep --image=nginx \$do > my-dep.yml` (deployment name is used as prefix for pods' name)

- Create a Service...to expose a...(faster than creating a service and editing its selector labels)
 - pod: `k expose pod my-pod --name my-svc [--type ClusterIp|NodePort|...] --port 3333 --target-port 80`
 - deployment: `k expose deployment nginx [--type ClusterIp|NodePort|...] --port=80 --target-port=8080`
- Note: A NodePort Service kind of lies on top of a ClusterIP one, making the ClusterIP Service reachable on the Node IPs (internal and external).
- Create Role / ClusterRole to permissions within a namespace / cluster-wide: `k create role my-role --verb=get,list,watch --resource=pods,pods/logs`
- Create RoleBinding / ClusterRoleBinding to connect Roles / ClusterRoles to subjects (users, groups or ServiceAccounts): `k create rolebinding my-rb --role=my-role --user=my-user`
- Create a service account, allowing container processes within Pods to authenticate with the K8s API: `k create sa my-sa`
- Create a quota: `k create quota my-q --hard=cpu=1,memory=1G,pods=2,services=3... [\$do]`

Update resources

- Add / delete / change a label: `k label (pods my-pod | nodes my-node) app=b` / `k label pods my-pod app=b` / `k label pods my-pod app=v2 --overwrite`
- Add a new label tier=web to all pods having 'app=v2' or 'app=v1' labels: `k label po -l "app in(v1,v2)" tier=web`
- Change a pod's image: `k set image my-pod nginx=nginx:1.7.1 [--record]`
- Recreate the pods in a deployment: `k rollout restart deploy my-dep`
- Perform a rolling update (e.g. to change an image): `k edit deployment my-dep` or `k set image deployment my-dep nginx=nginx:1.21.5 [--record]`
- Check rollout status: `k rollout status deployment my-dep`
- Roll back to the previous version: `k rollout undo deployment my-dep`
- Scale a deployment [and record the command (into Annotations > change-cause)]: `k scale deployment my-dep --replicas=5 [--record]`
- Autoscale a deployment, from 5 to 10 pods, targeting CPU utilization at 80%: `k autoscale deploy my-dep --min=5 --max=10 --cpu-percent=80`
 - View the Horizontal Pod Autoscalers (hpa): `k get hpa nginx`

Debugging

- Use `k get pods [-A] [--show-labels]`: check `STATUS`, `READY` and `RESTARTS` attributes.
- Retrieve a pod status: `k get pod <pod_name> -o json | jq .status.phase`
- Retrieve pod names, their resources and QoS class: `k get pod -o jsonpath="{range .items[*]}{.metadata.name} - {.spec.containers[*].resources} - {.status.qosClass}{'\n'}}`
- Retrieve pod / container logs: `k logs <pod_name> [-c <container_name>] [-p]` (if pod crashed and restarted, -p option gets logs about the previous instance)
- List events for a / all namespace(s), sorted by time: `k get events (-n <my-namespace> | -A) [--sort-by=.metadata.creationTimestamp]`
- Show metrics for pods (including containers) / nodes: `k top pod [--containers] [--sort-by (cpu | memory)] [-l app=b]` / `k top node [--sort-by (cpu | memory)]`
- Warning `Back off restarting failed container x in pod y`: possible cause is a missing `command` for a busybox container

Delete / replace resources

- Force replace a resource: `k replace --force -f ./pod.json`
- Delete pods and services using their label: `k delete pods,services -l app=b \$now`

Secrets

- Use `k get secret ...` to get a base64 encoded token
- Use `k describe secret ...` to get a base64 decoded token...or pipe it manually through `echo <token> | base64 -d -`
- If a Secret belongs to a ServiceAccount, it'll have the annotation `kubernetes.io/service-account.name`

Networking, services, DNS

- The cluster has a single virtual network spanning across all nodes.
- Nodes remain `NotReady`, unable to run Pods, until a network plugin is installed. `Starting kube-proxy` is shown in the nodes logs and no networking pods exist.
- Default FQDN:
 - `.<my-namespace>.pod.cluster.local.`
 - `.<my-namespace>.svc.cluster.local.`
- `from` and `to` selectors:
 - ![(np_from_to_selectors.png)]
- Example of "deny all" policy for labelled pods:


```
...
spec:
  podSelector:
    matchLabels:
      app: maintenance
  policyTypes:
    - Ingress
    - Egress
  ...
```

```
spec:
  podSelector:
    matchLabels:
      app: maintenance
  policyTypes:
    - Ingress
    - Egress
  ...
```

- Example of egress policy, 1) restricting outgoing tcp connections from frontend to api, 2) still allowing outgoing traffic on UDP/TCP ports 53 for DNS resolution.

```
spec:
  podSelector:
    matchLabels:
      id: frontend          # label of the pods this policy should be applied on
  policyTypes:
    - Egress                # we only want to control egress
  egress:
    - to:                   # 1st egress rule
      - podSelector:        # allow egress only to pods with api label
```

```

    matchLabels:
      id: api
- ports:
- port: 53
  protocol: UDP
- port: 53
  protocol: TCP
  # 2nd egress rule
  # allow DNS UDP
  # allow DNS TCP
...
- Example policy allowing all pods in the `users-backend` namespace to communicate with each other only on a
specific port (80): first label the namespace: `k label namespace users-backend app=users-backend`:
...
metadata:
  name: np-users-backend-80
  namespace: users-backend
spec:
  podSelector: {}
  policyTypes:
  - Ingress
  ingress:
  - from:
    - namespaceSelector:
        matchLabels:
          app: users-backend
    ports:
    - protocol: TCP
      port: 80
  # selects all pods in the specified namespace
...
- Example policy containing a single `from` element allowing connections from Pods with the label `role=
client` in namespaces with the label `user=alice`
...
  ingress:
  - from:
    - namespaceSelector:
        matchLabels:
          user: alice
    podSelector:
      matchLabels:
        role: client
  ...
- Example policy containing two elements in the `from` array, allowing connections from Pods in the local
Namespace with the label `role=client`, **or** from any Pod in any namespace with the label `user=alice`.
...
  ingress:
  - from:
    - namespaceSelector:
        matchLabels:
          user: alice
    - podSelector:
        matchLabels:
          role: client
  ...
- When in doubt, `kubectl describe` shows how Kubernetes has interpreted the policy.
- Endpoints are the underlying entities (such as Pods) that a Service routes traffic to.
- Ingress: manages external access to Services; more powerful than a simple NodePort Service (e.g. SSL
termination, advanced load balancing, or namebased virtual hosting).

### Helm
- List release with `helm [-n my_ns] ls [-a]`
- List pending deployments on all namespaces: `helm list --pending -A`
- List / search repo: `helm repo list` / `helm search repo nginx`
- Download (not install) a chart from a repository: `helm pull [chart URL | repo/chartname] [...] [flags]`
- Untar a chart (after downloading it): `helm pull --untar [rep/chartname]`
- Check customisable values setting for an install, e.g. `helm show values bitnami/apache [| yq e]`
- Custom install example `helm install my-apache bitnami/apache --set replicaCount=2`
- Upgrade a release, e.g. `helm upgrade my-api-v2 bitnami/nginx`
- Undo a helm rollout/upgrade: `helm rollback`
- Delete an installed release with `helm uninstall <release_name>`

### etcd
- etcd is a consistent and highly-available key value store used to store for all cluster data.
- Backup / restore etcd data: [link](CKA%20training/etcd.md)

```