



KubeCon



CloudNativeCon

North America 2019

PodOverhead: Accounting for Greater Cluster Stability

Eric Ernst, Intel, @egernst

SW Engineer
kat herder
kata containers arch committee
k8s contributor



What happens when you click enter after typing `kubectl apply` for a pod?

Where and what are these overheads?

What is the PodOverhead feature?

Why should you care?

CLIENT

my-yaml
podSpec

KUBECTL

CLIENT

my-yaml
podSpec

KUBECTL

- (1) verify yaml
- (2) create request

CLIENT

my-yaml
podSpec

KUBECTL

- (1) verify yaml
- (2) create request

HTTP request

MASTER NODE

API SERVER

CONTROLLER
MANAGER

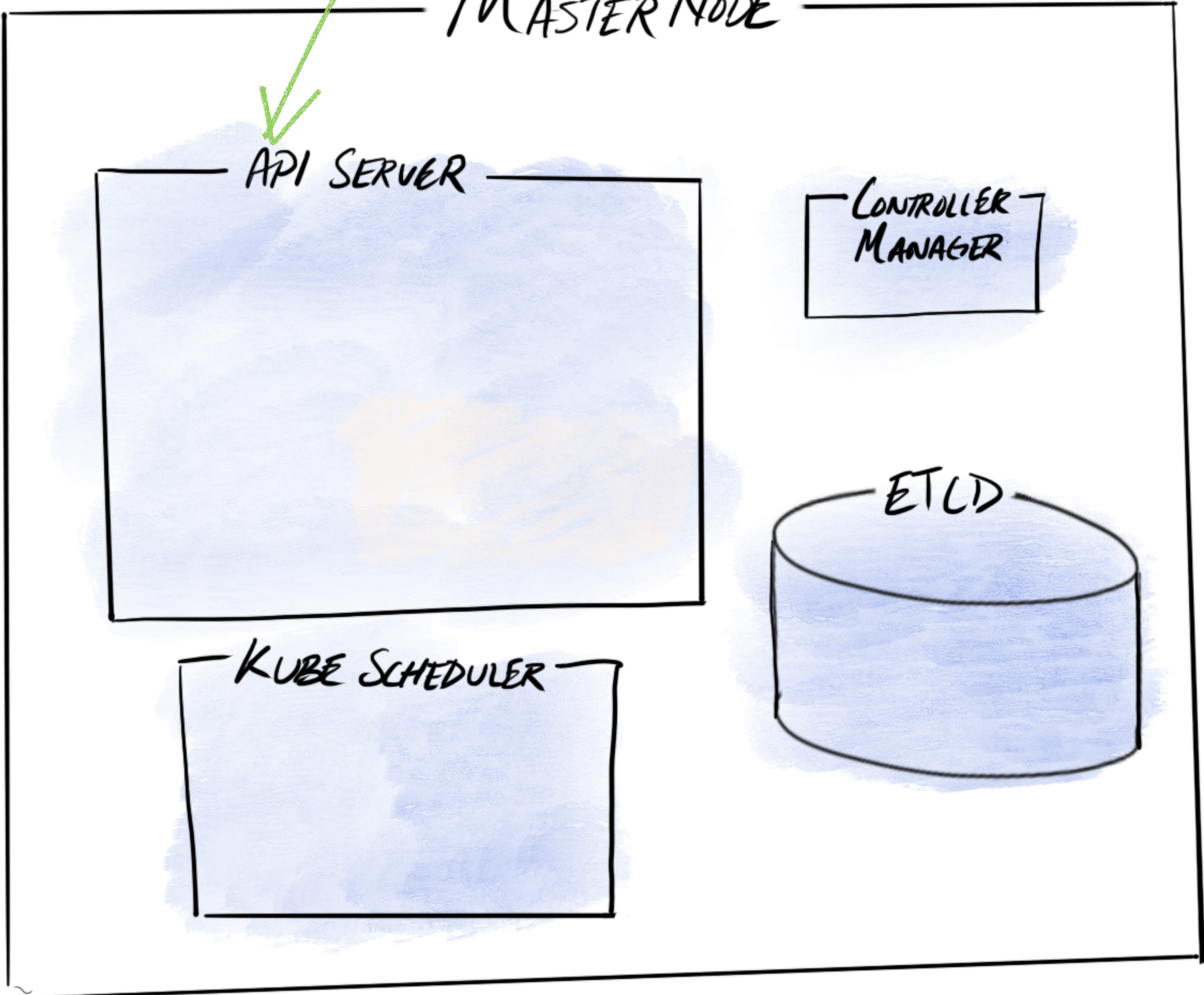
MASTER NODE

API SERVER

CONTROLLER
MANAGER

ETCD

KUBE SCHEDULER



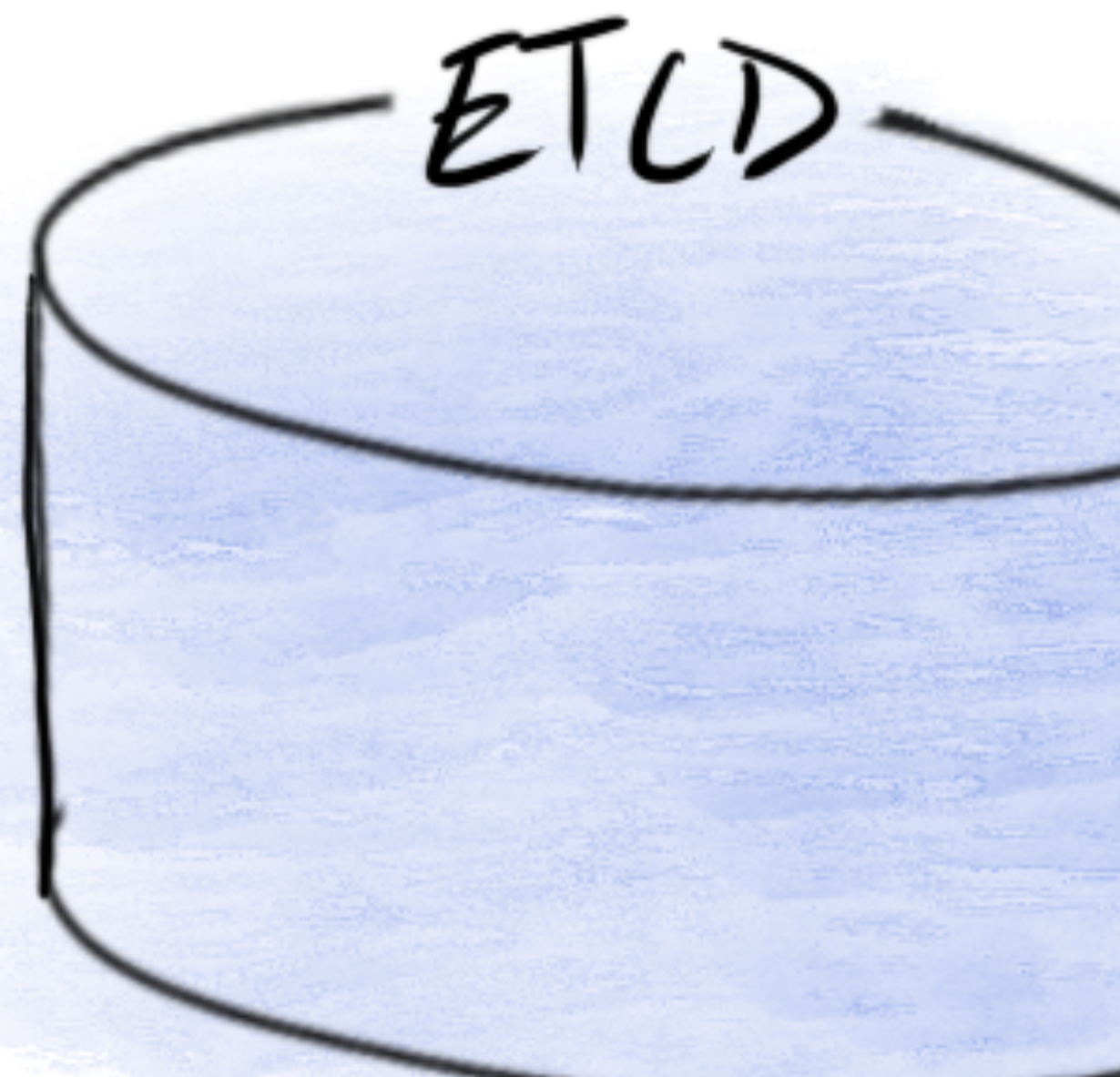
MASTER NODE



API SERVER



CONTROLLER
MANAGER

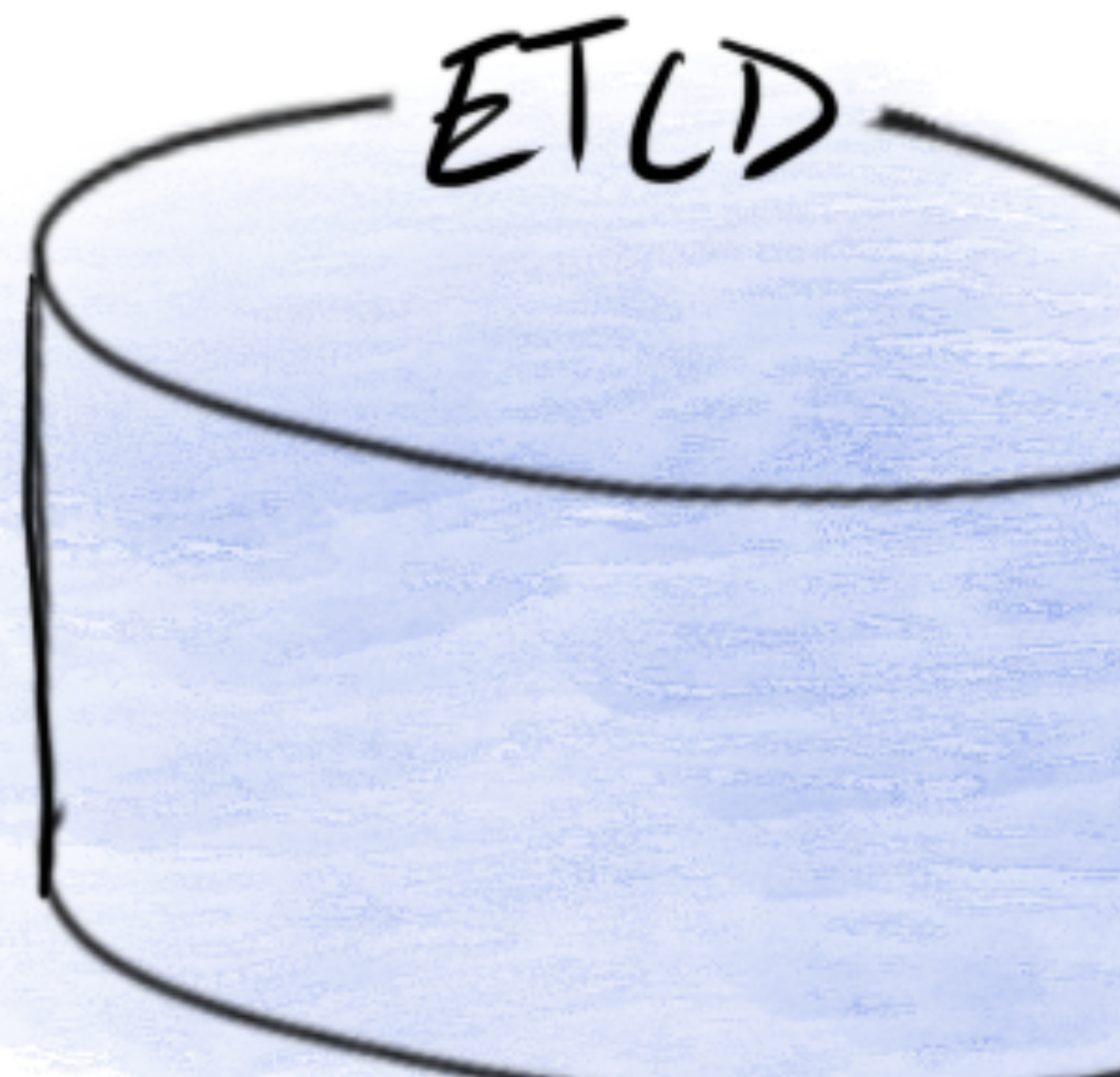
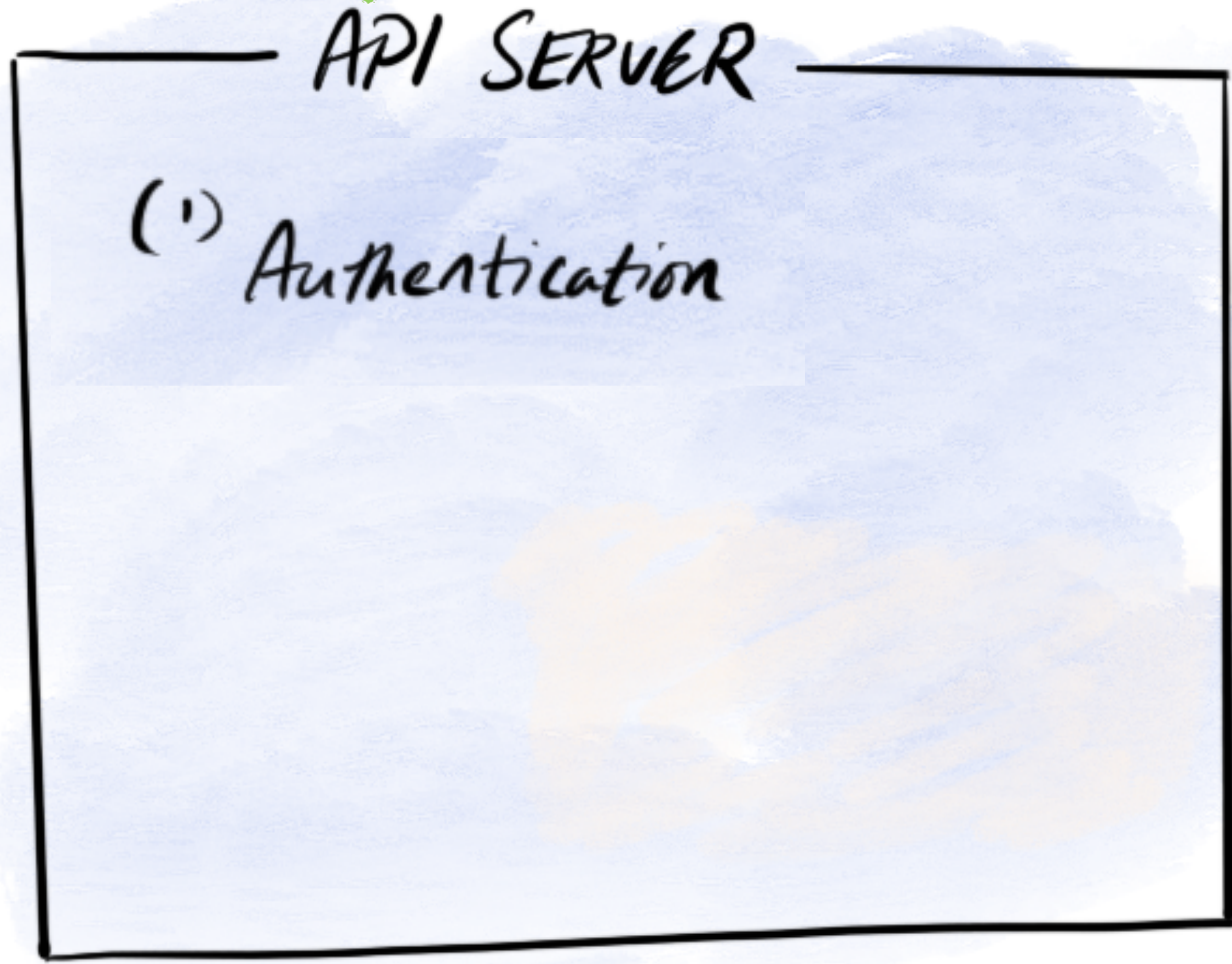


ETCD

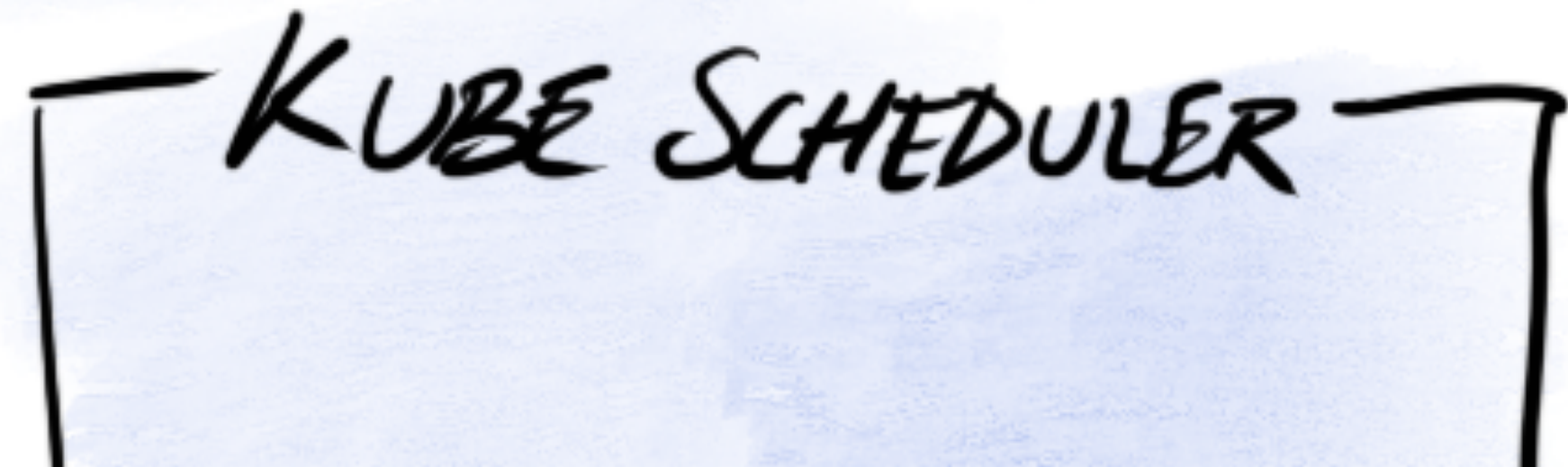
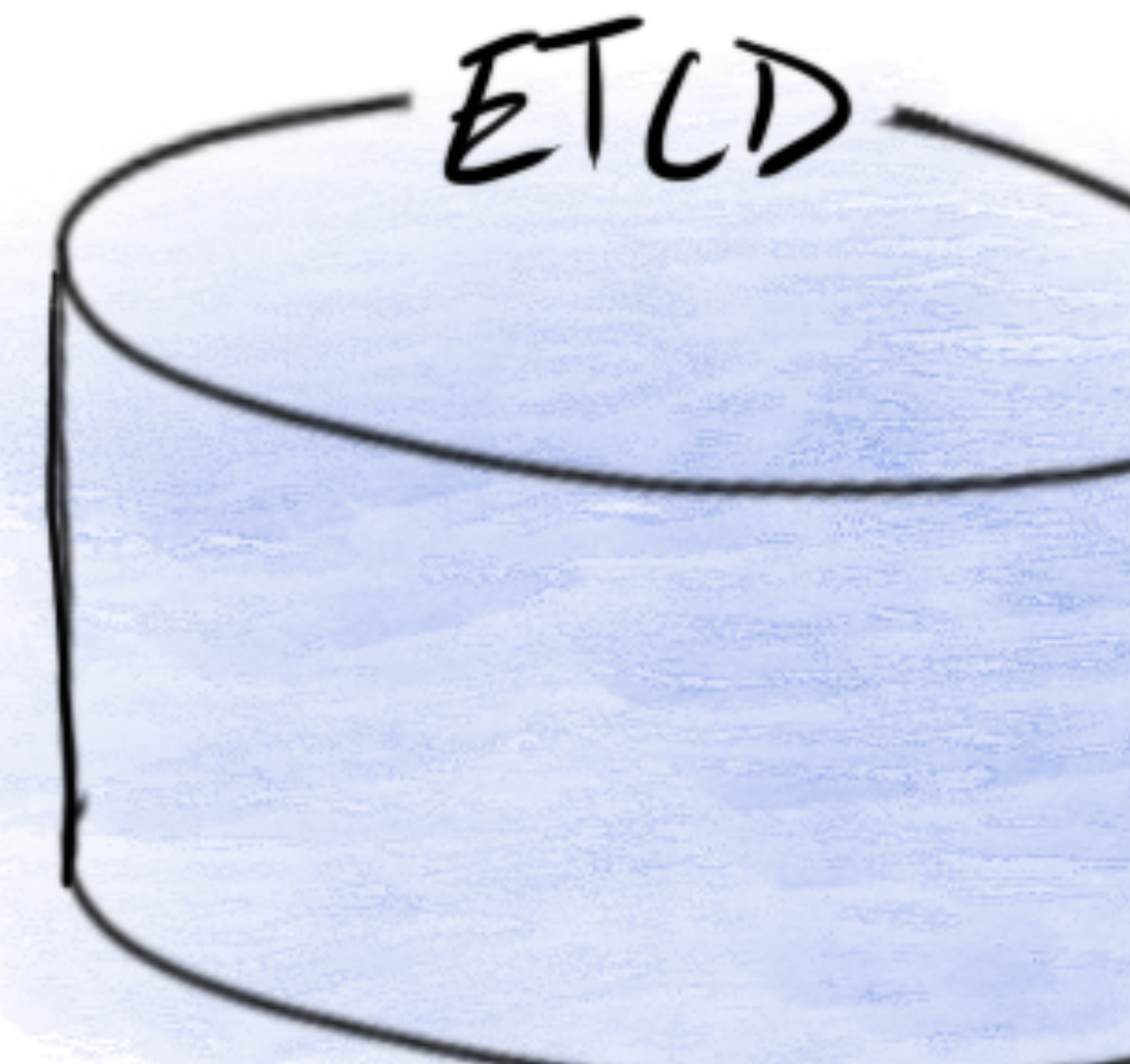
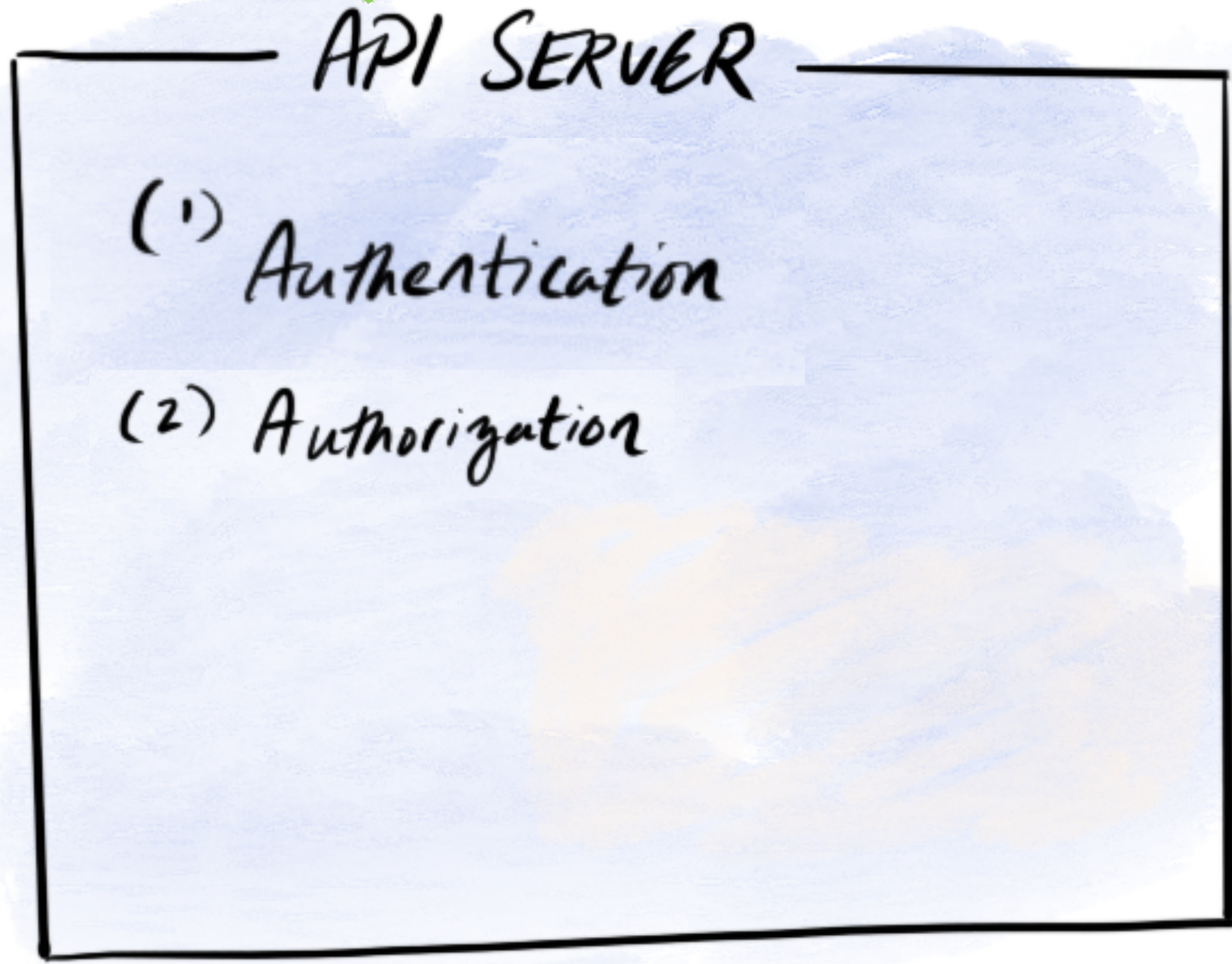


KUBE SCHEDULER

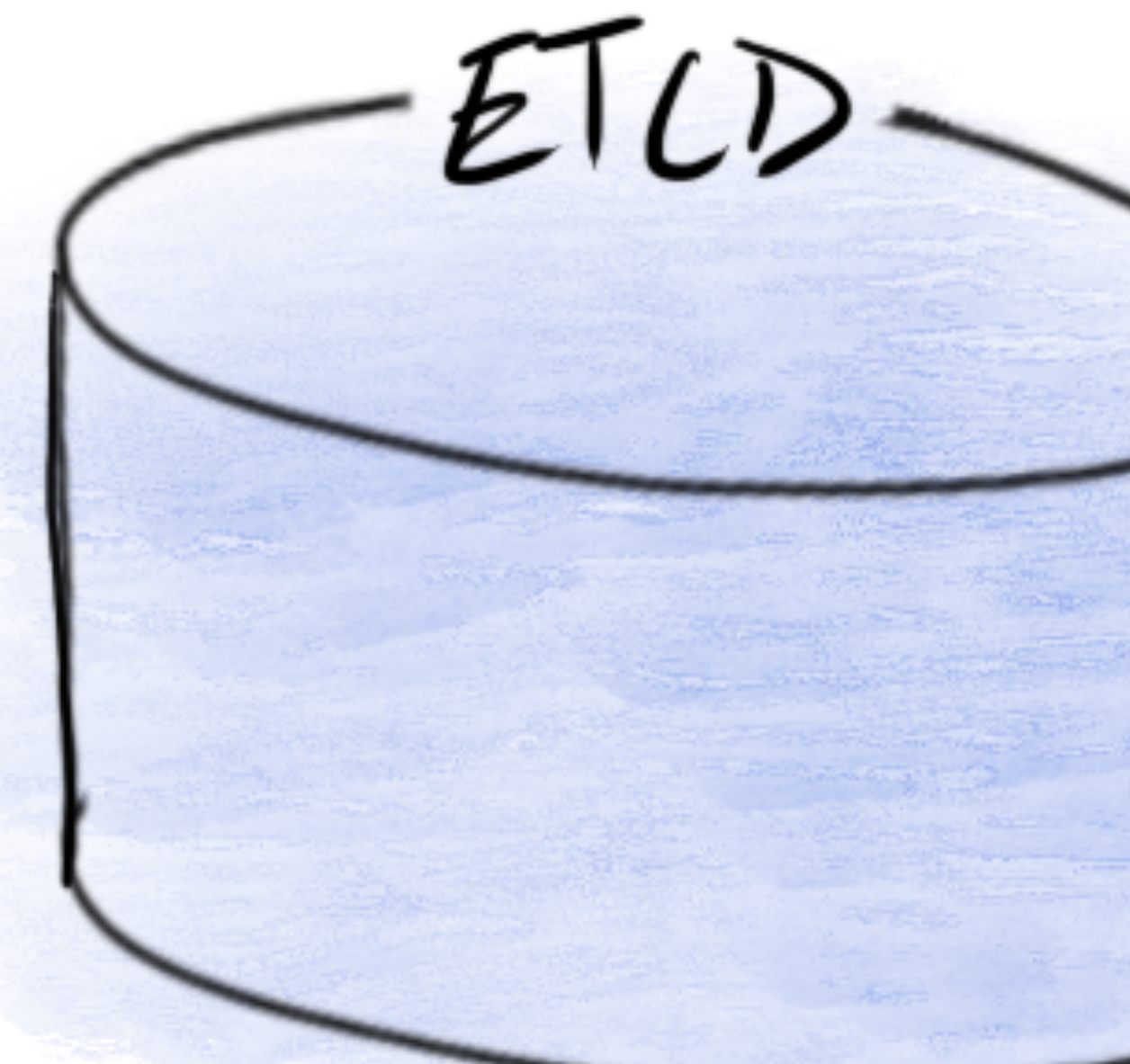
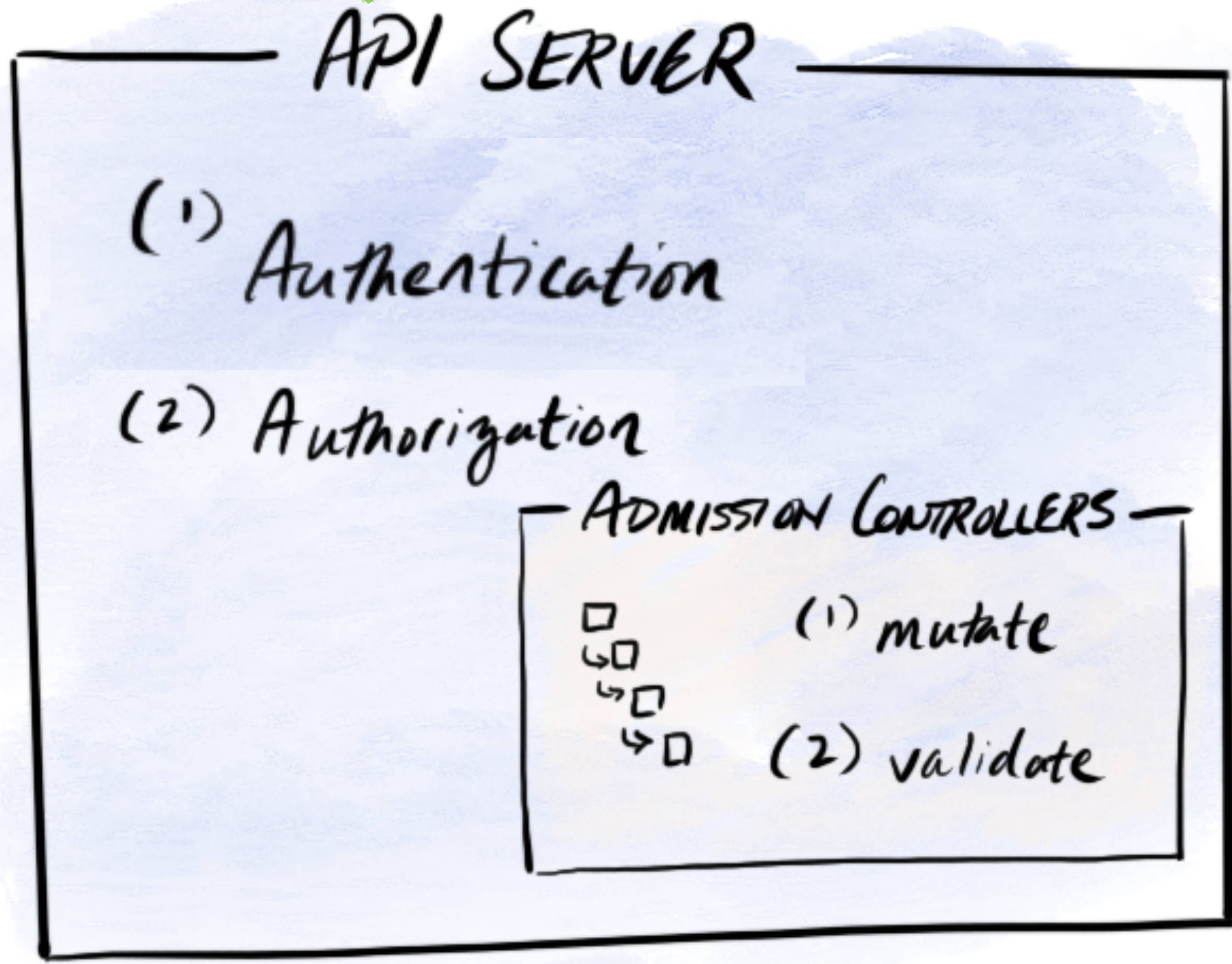
MASTER NODE



MASTER NODE

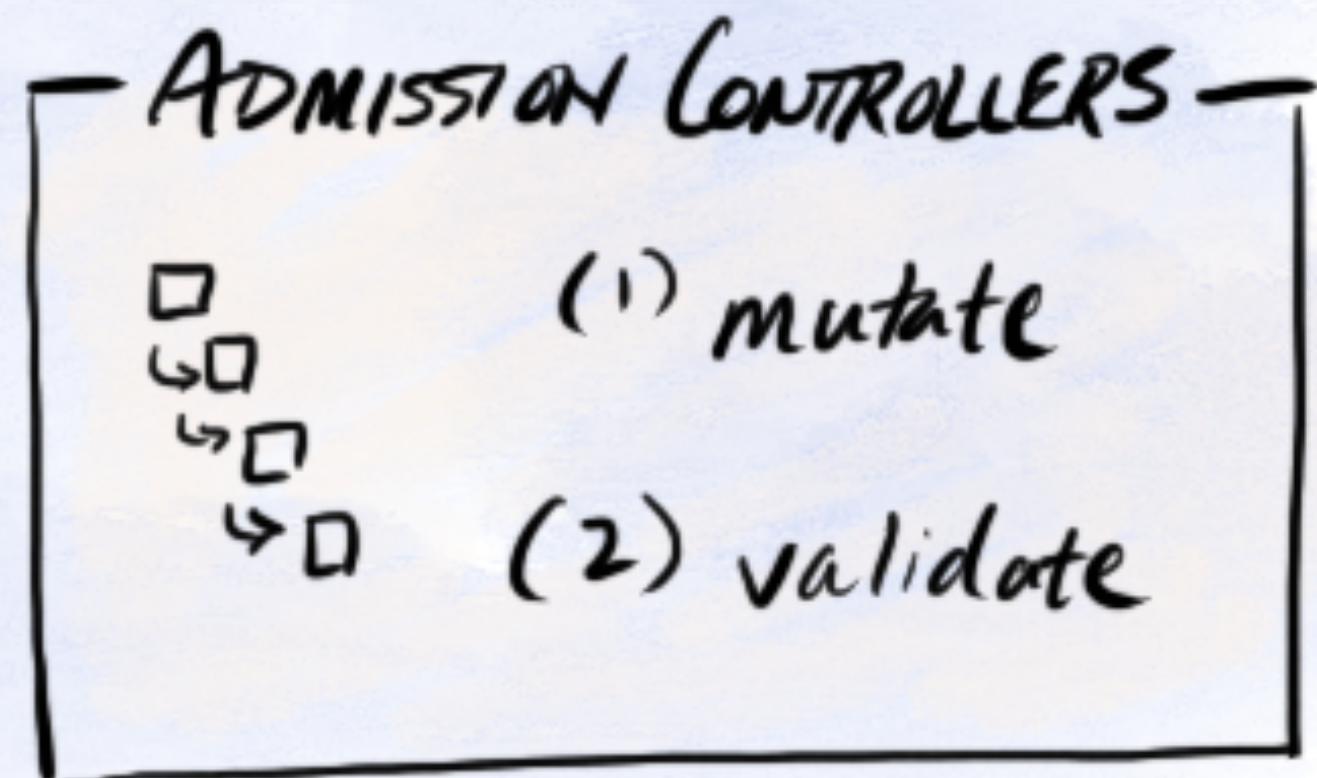


MASTER NODE



Authentication

Authorization



MANAGER

ETCD

Pod pod-list
: nodeName
awesome :

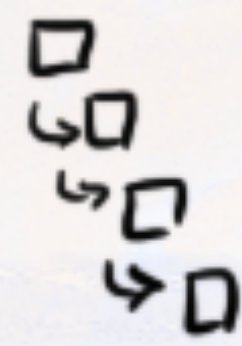
KUBE SCHEDULER

```
kubectl get pods --all-namespaces
```

NAMESPACE	NAME	READY	STATUS	RESTARTS	AGE
default	awesome-pod	0/1	Pending	0	100 years

(2) Authorization

ADMISSION CONTROLLERS



(1) mutate

(2) validate

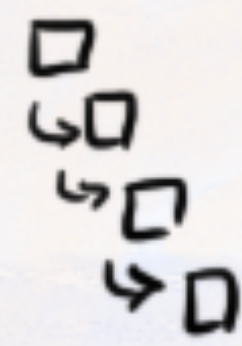
KUBE SCHEDULER

ETCD

Pod-list
Pod : nodeName
awesome :

(2) Authorization

ADMISSION CONTROLLERS



(1) mutate

(2) validate

KUBE SCHEDULER

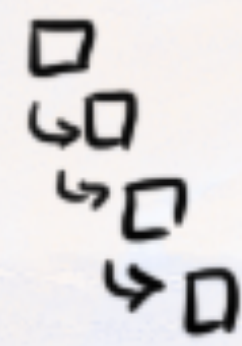
(1) predicate }
(2) prioritize } ↻

ETCD

Pod-list
Pod : nodeName
awesome :

(2) Authorization

ADMISSION CONTROLLERS



(1) mutate

(2) validate

KUBE SCHEDULER

(1) predicate }
(2) prioritize } ↻

ETCD

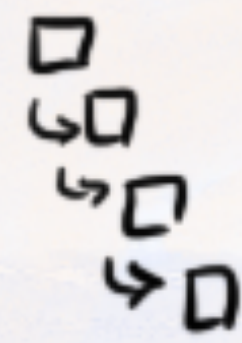
Pod-list
Pod : nodeName
awesome : worker

```
kubectl get pods --all-namespaces
```

NAMESPACE	NAME	READY	STATUS	RESTARTS	AGE
default	awesome-pod	0/1	PodScheduled	0	100 years

(2) Authorization

ADMISSION CONTROLLERS



(1) mutate

(2) validate

KUBE SCHEDULER

(1) predicate }
(2) prioritize } ↻

ETCD

Pod-list
Pod : nodeName
awesome : worker

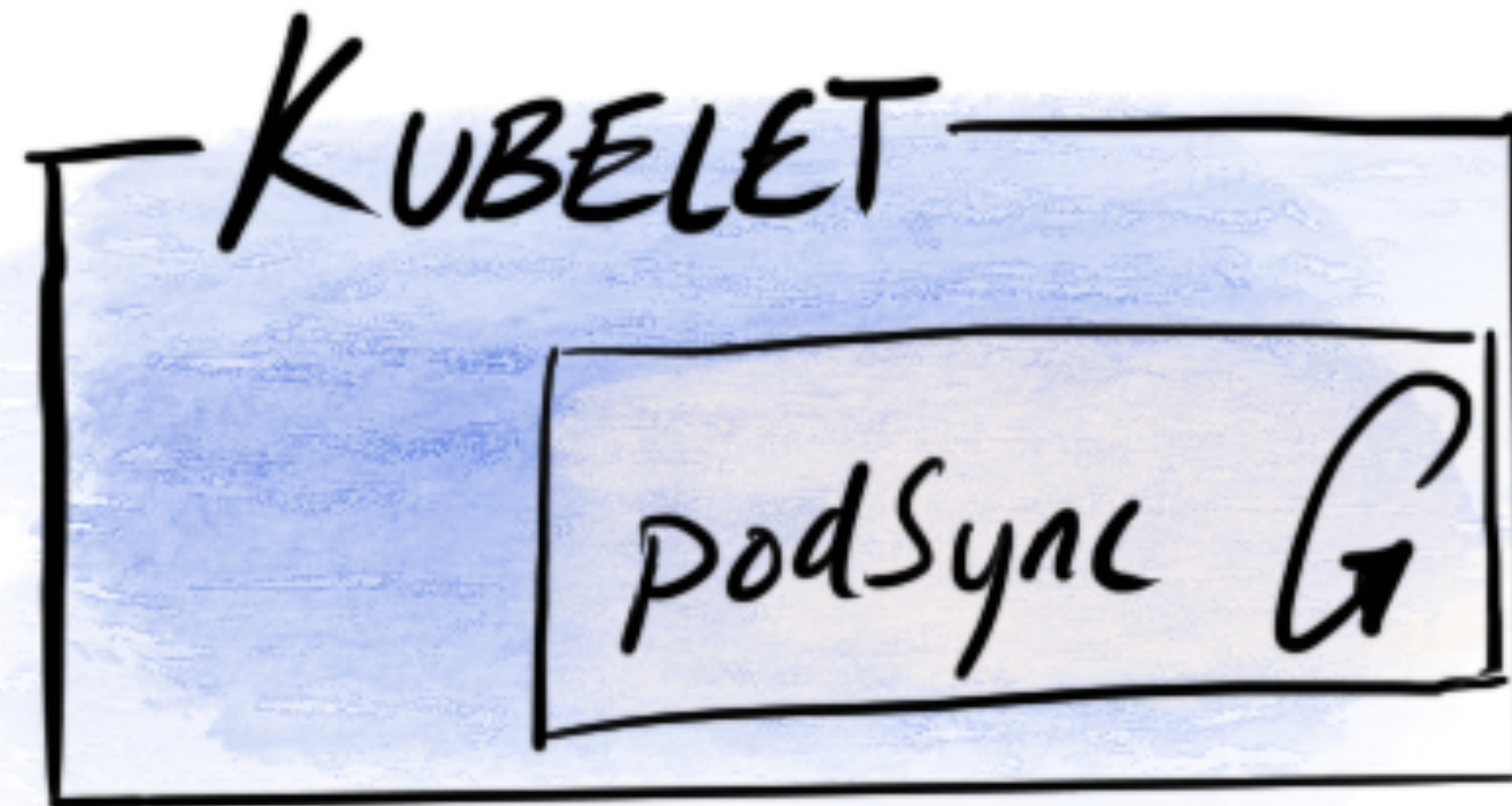
WORKER NODE

KUBELET

podSync G

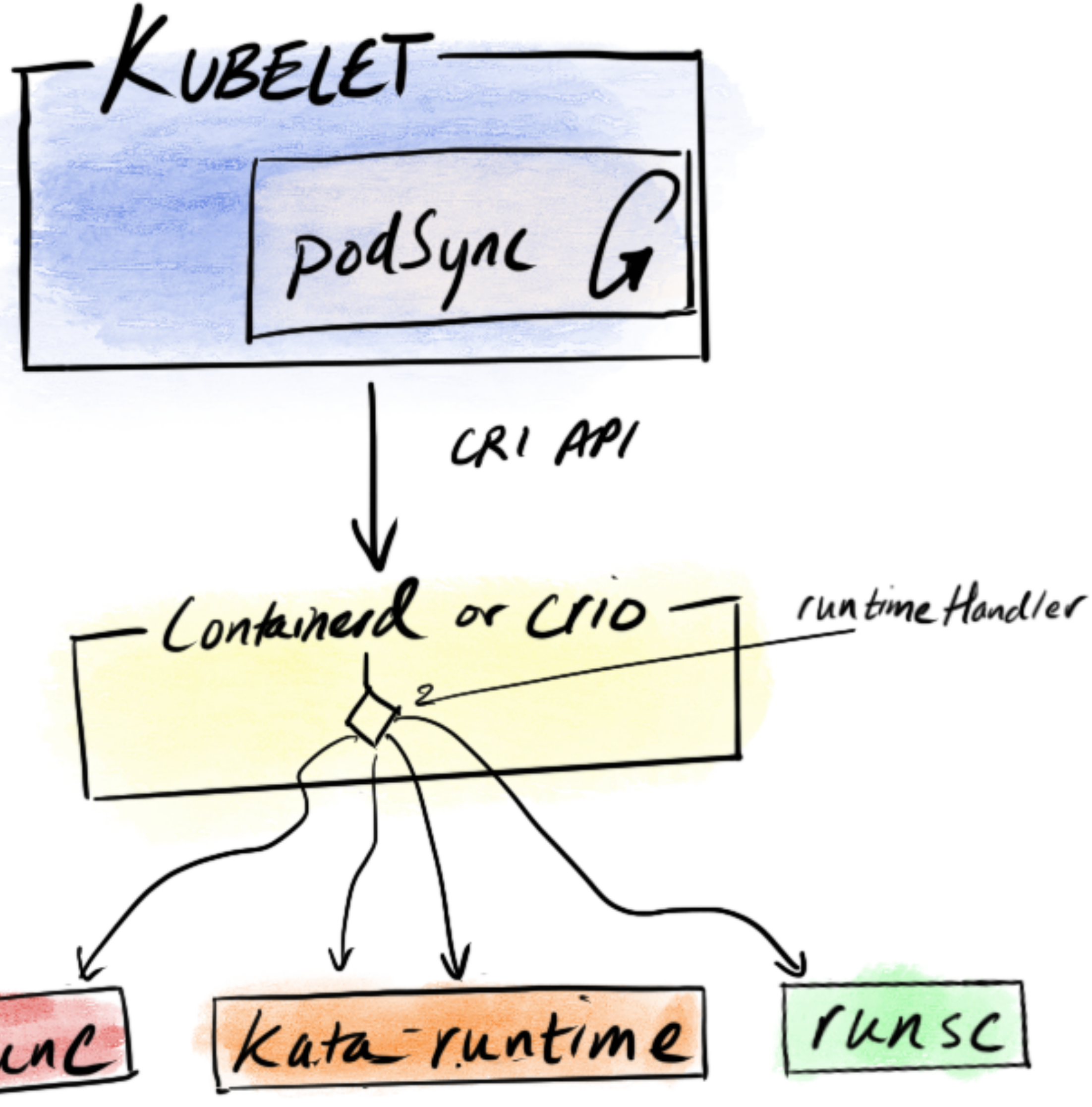
CRI API

WORKER NODE



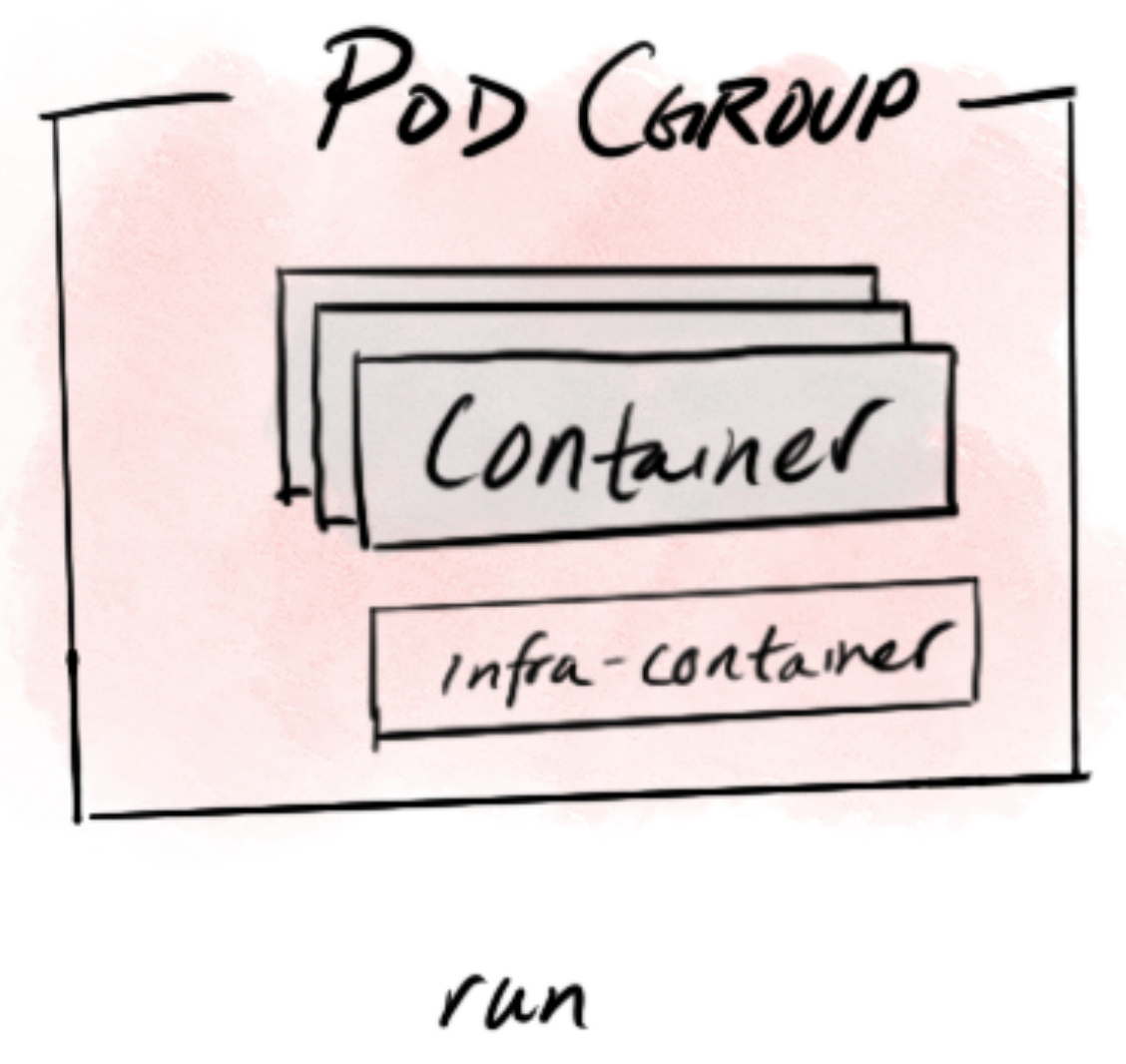
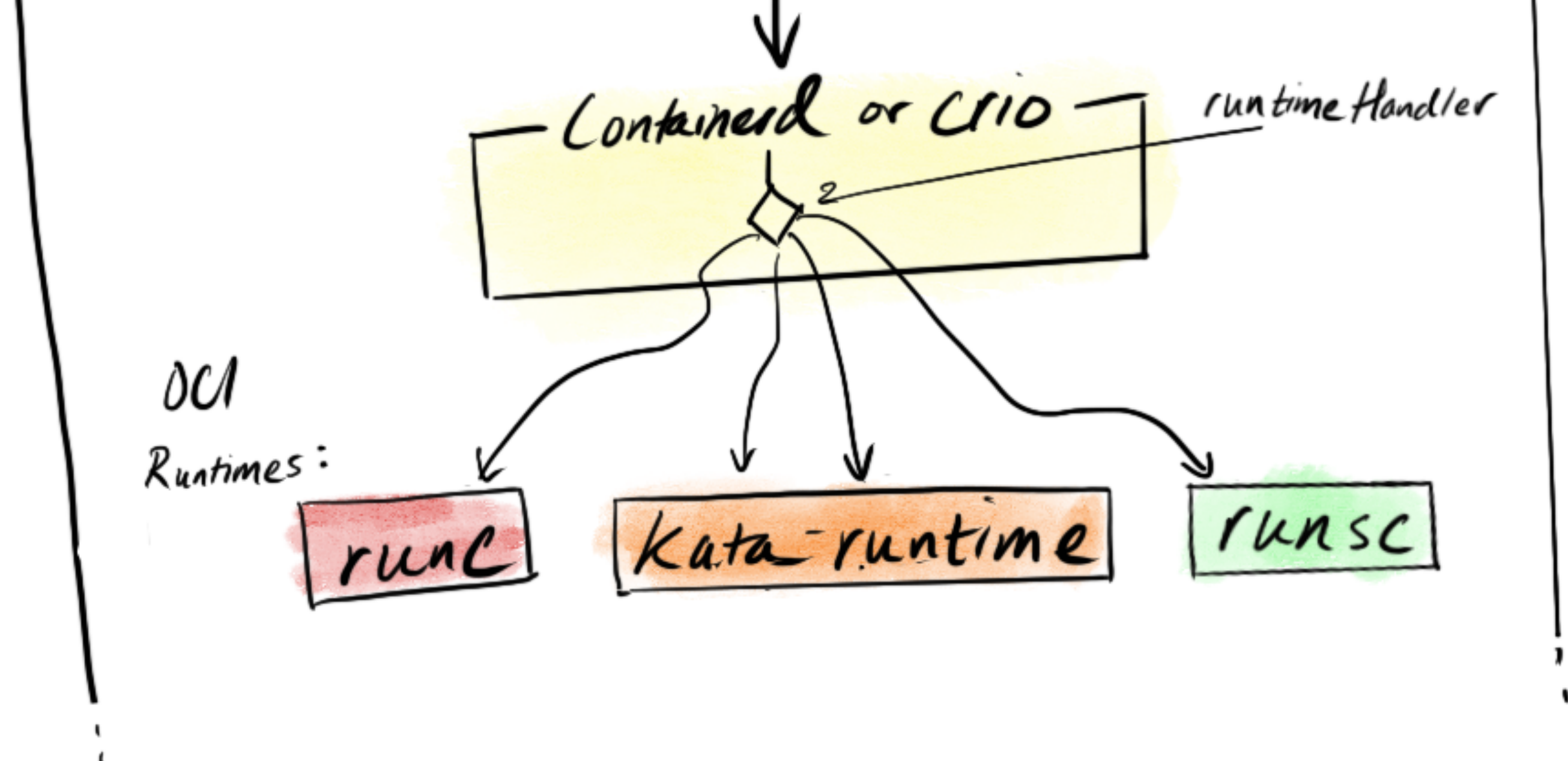
CRI API

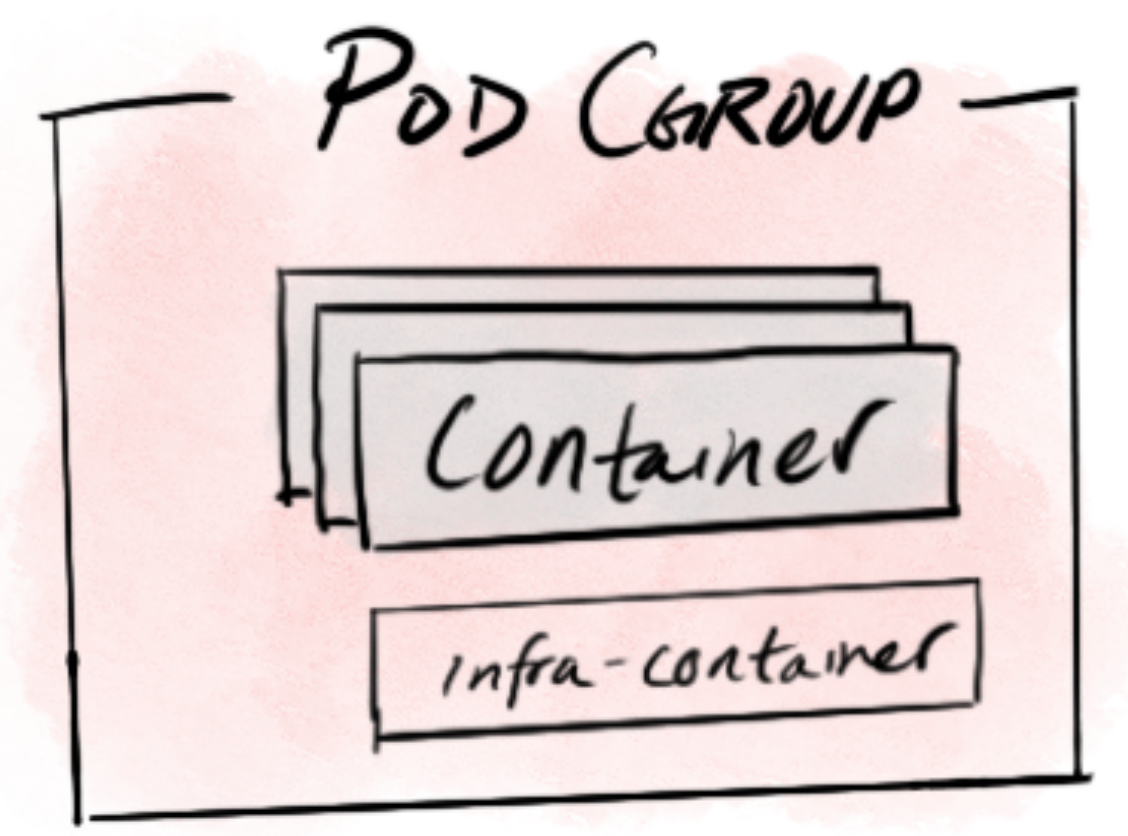
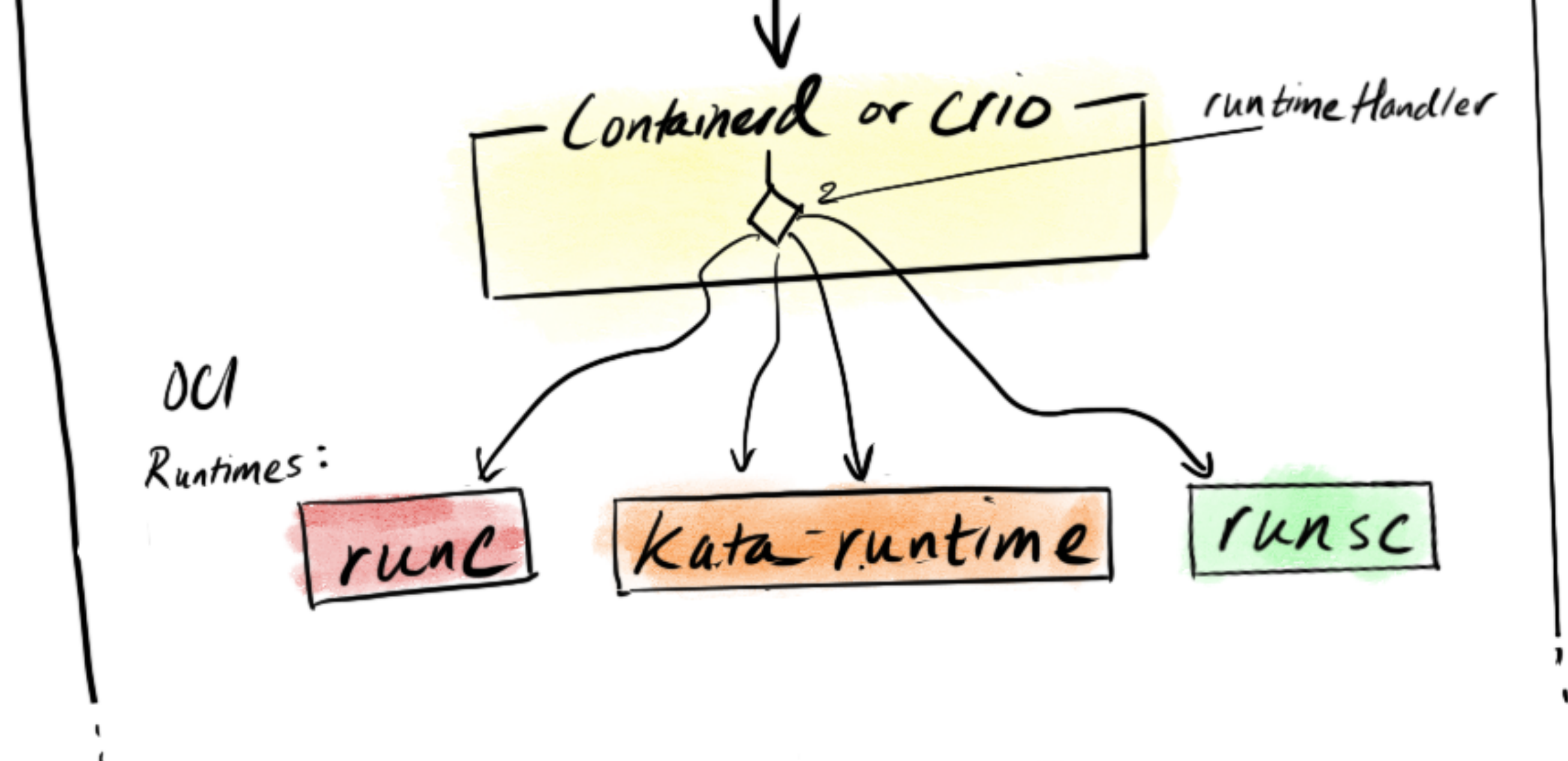




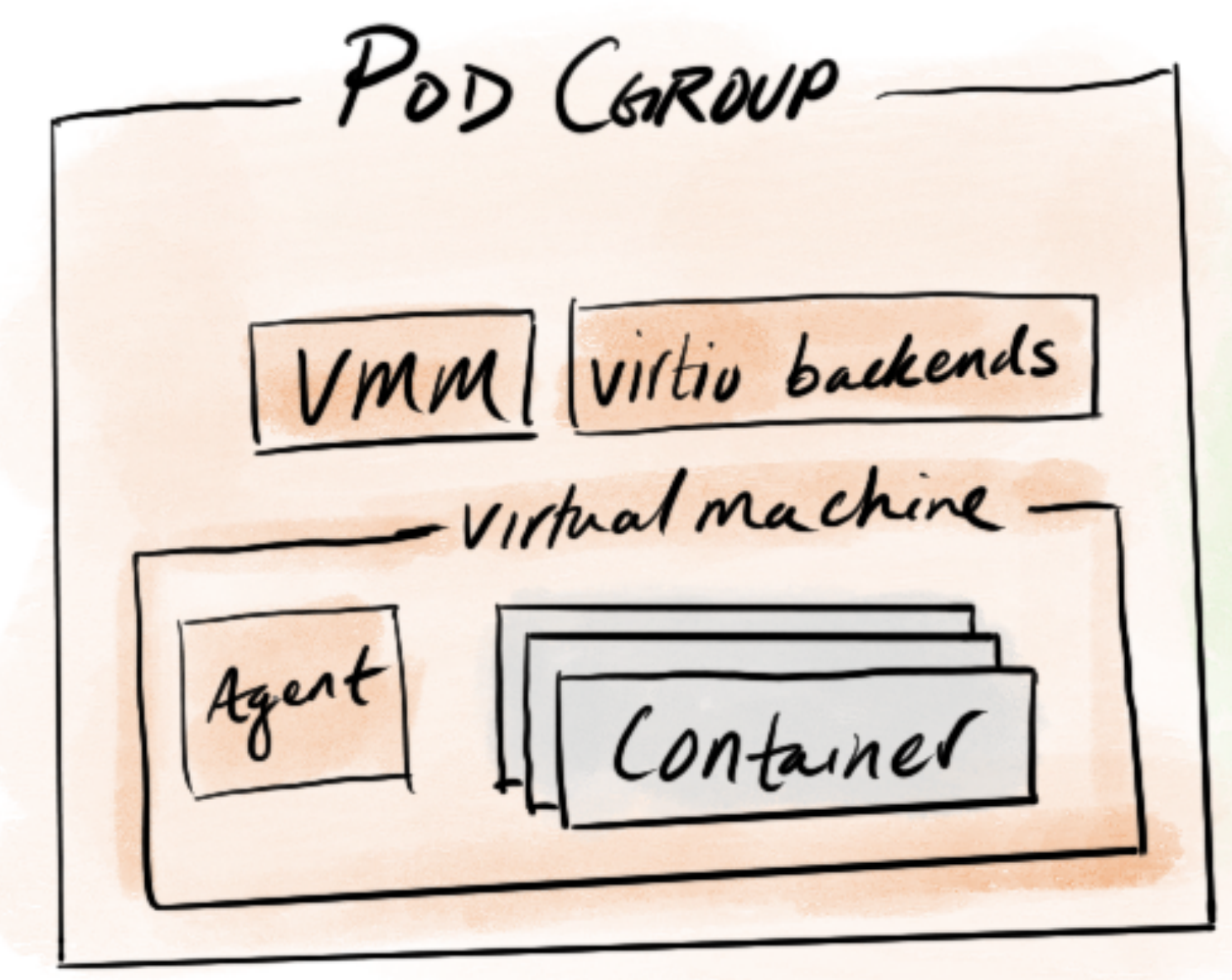
- POD SYNC -

- 1) Create Pod cgroup
- 2) Ask CRI to create sandbox
 - ↳ CRI calls CNI
 - ↳ CRI calls OCI
- 3) Ask CRI to pull image
- 4) Ask CRI to create container
 - ↳ CRI calls OCI

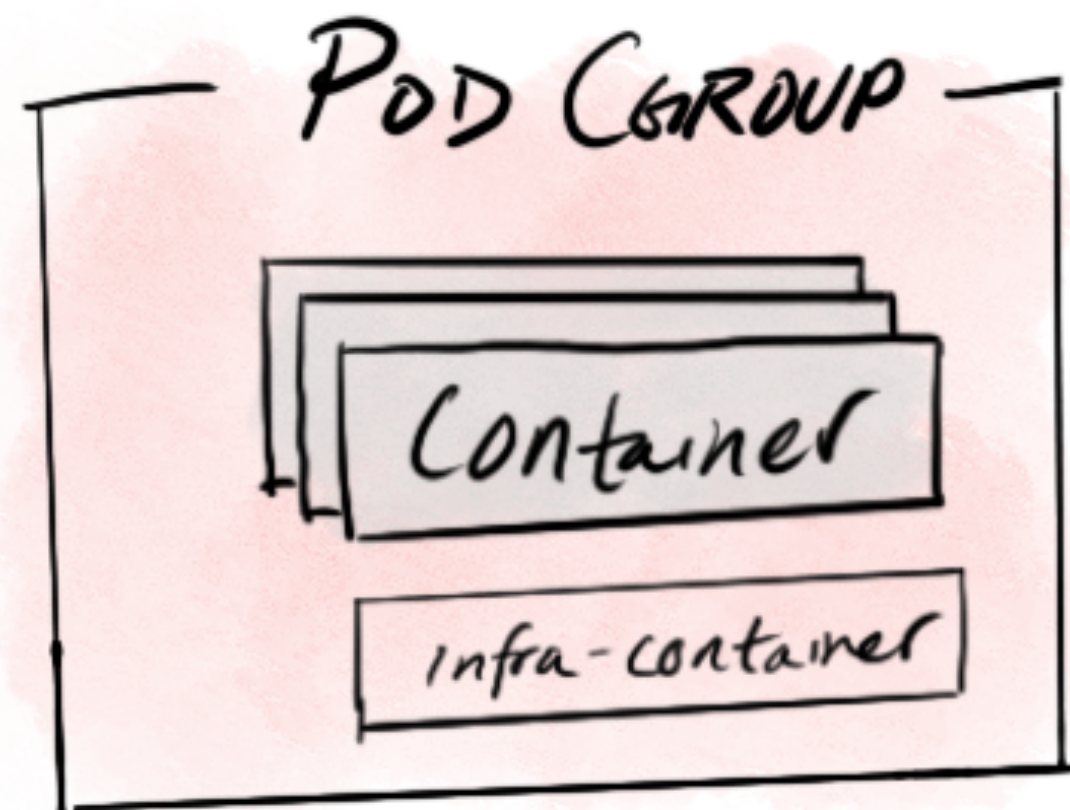
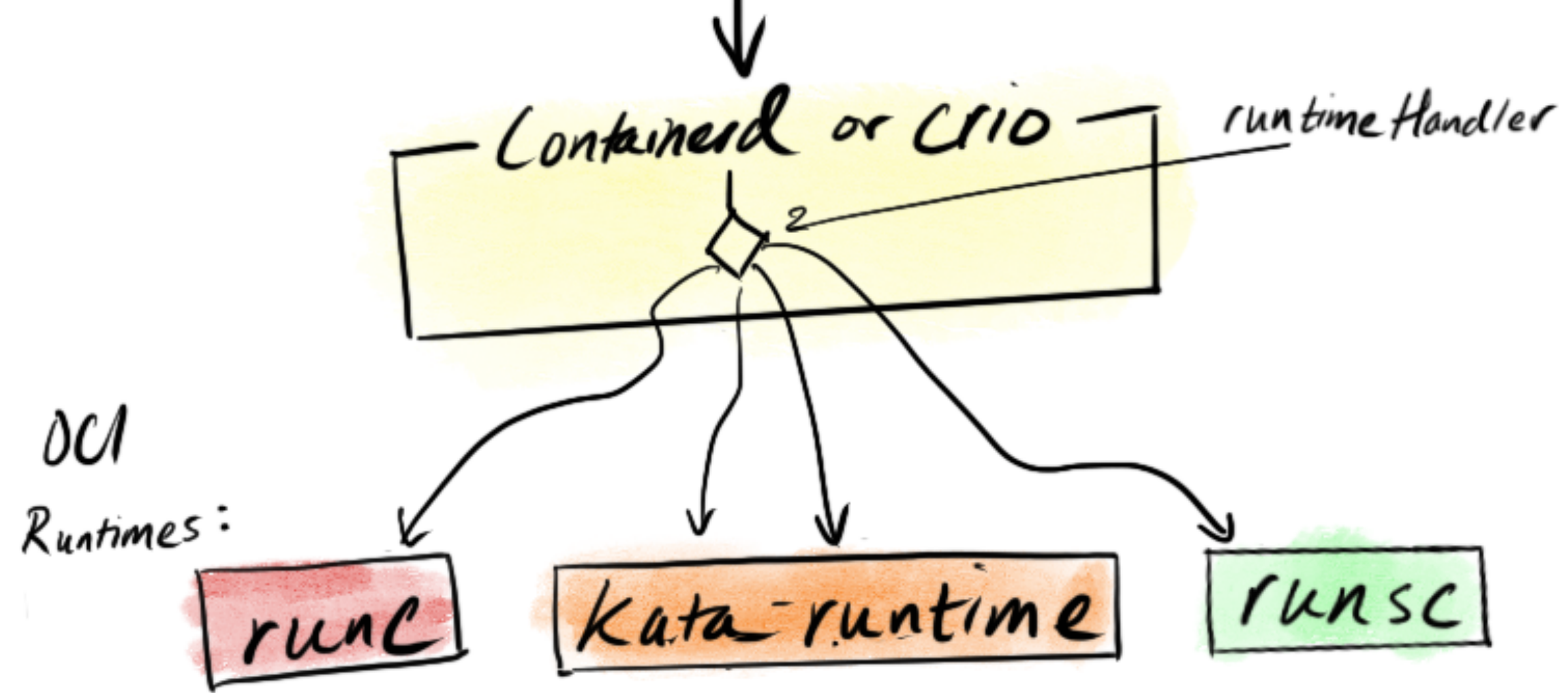




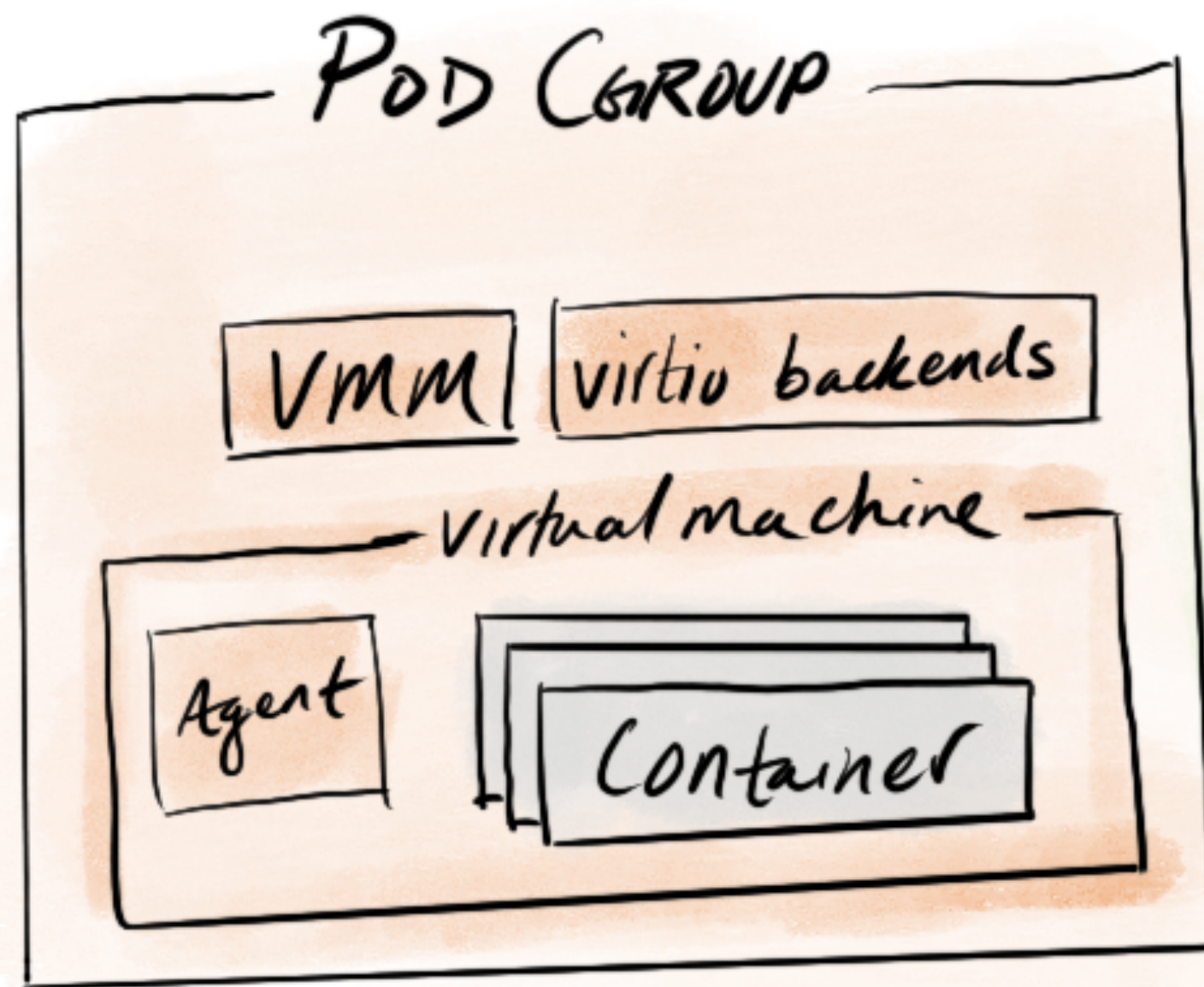
run



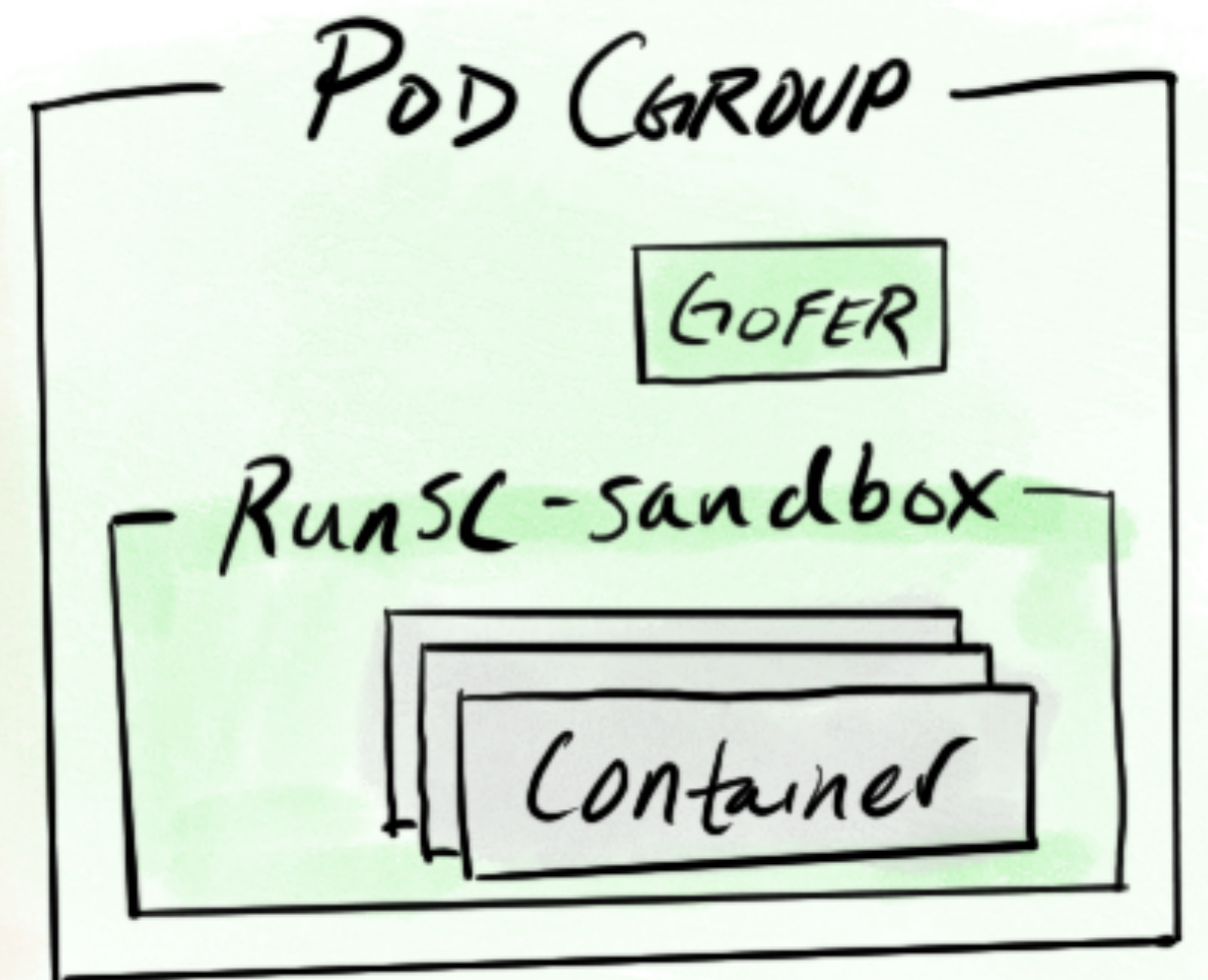
Kata



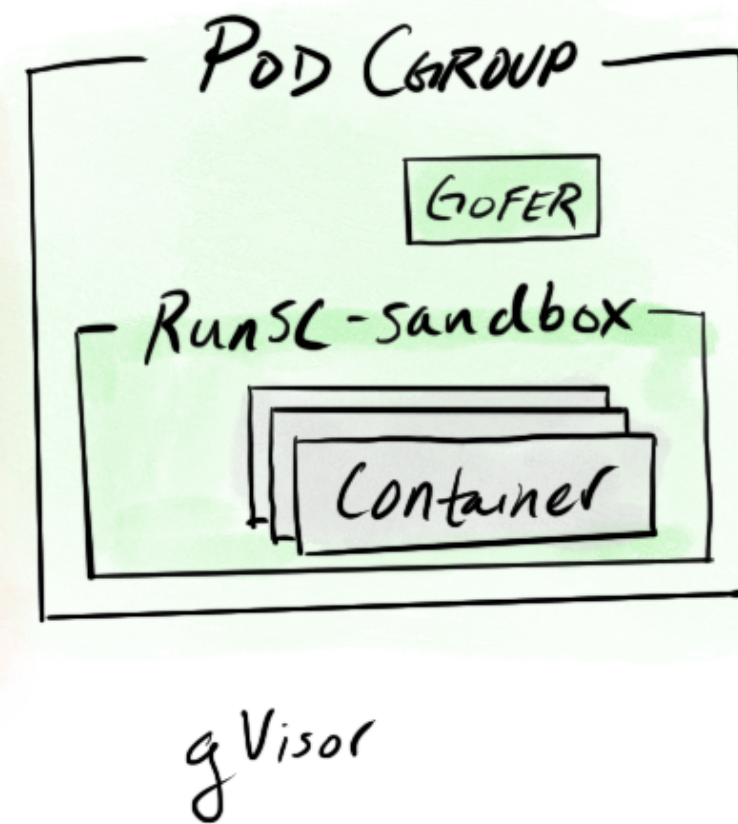
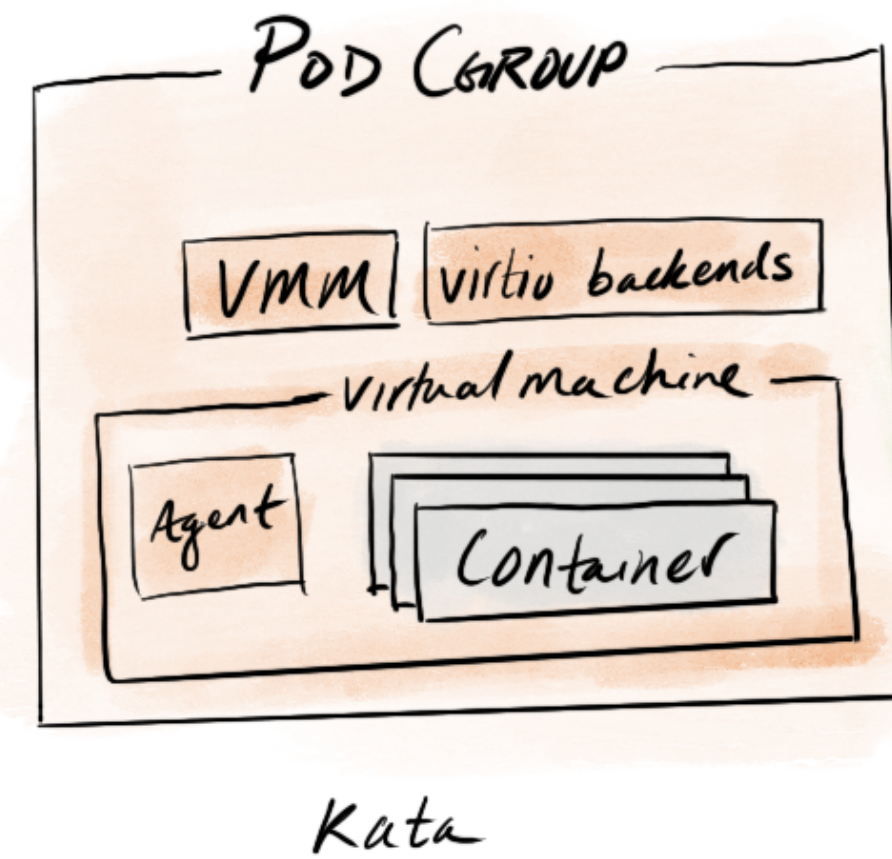
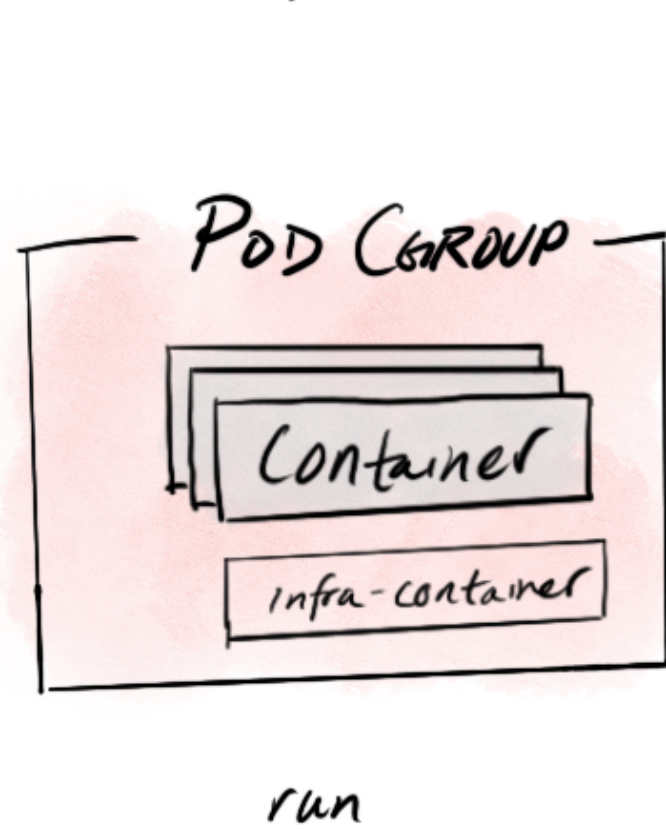
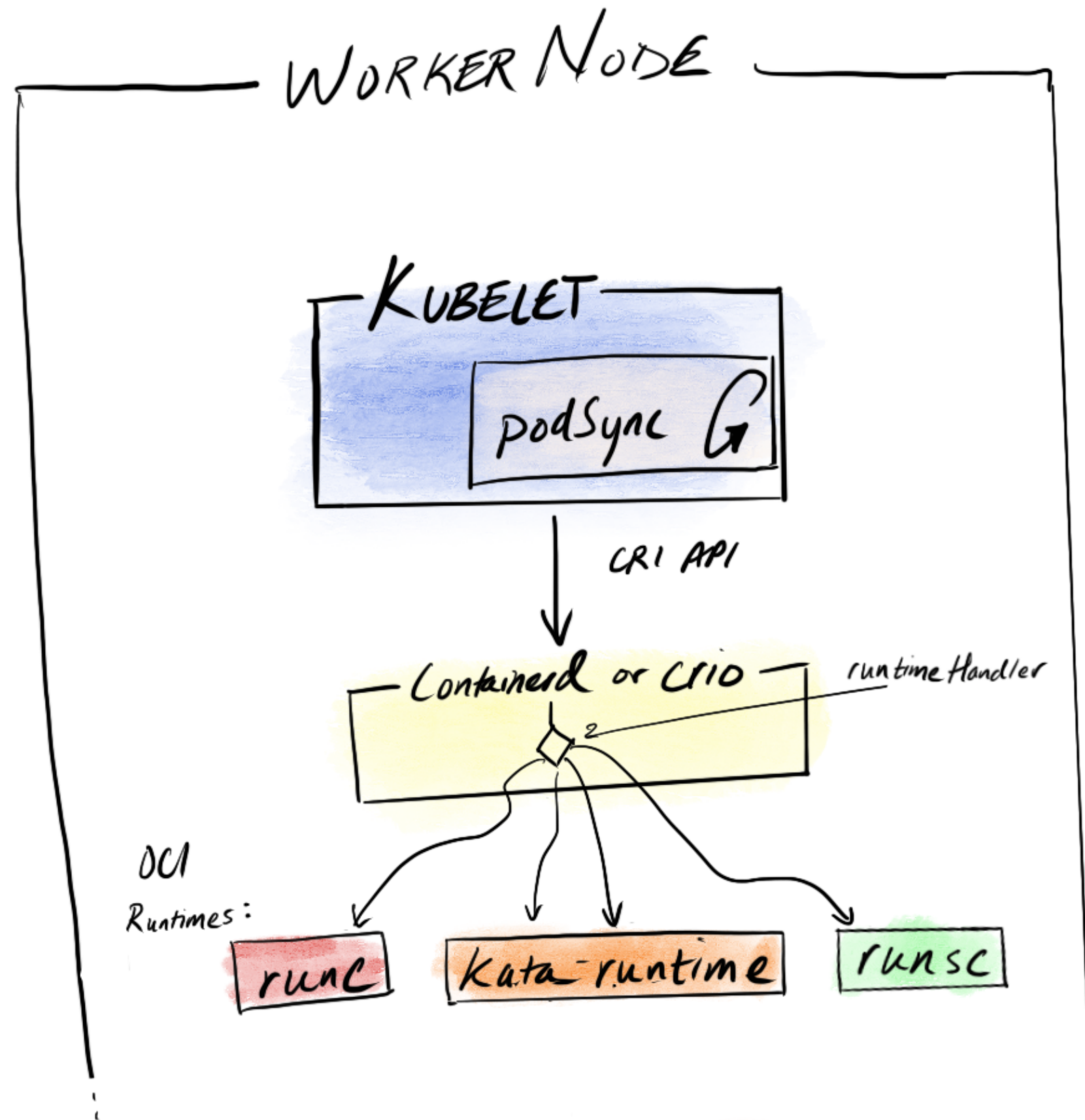
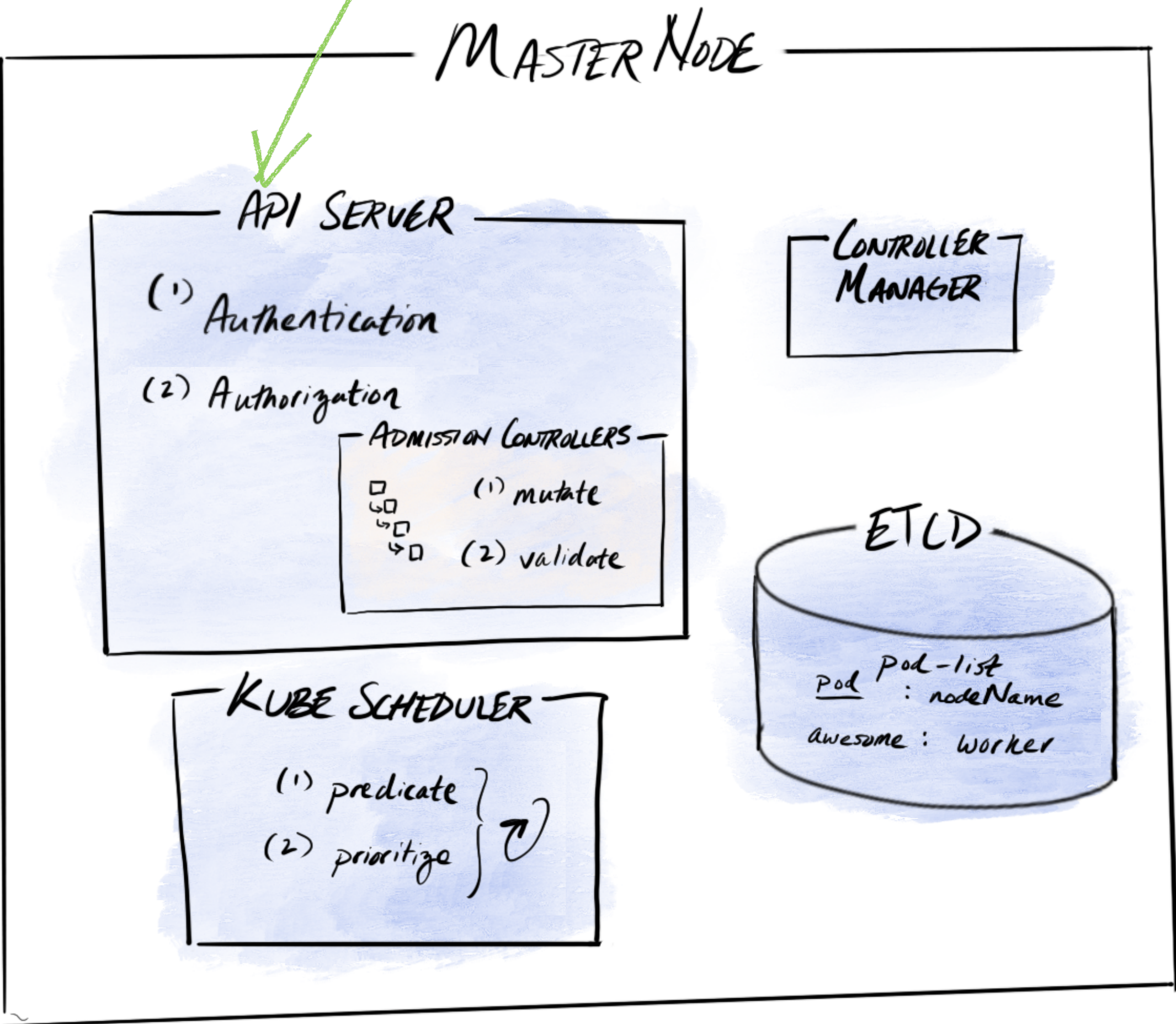
run



Kata

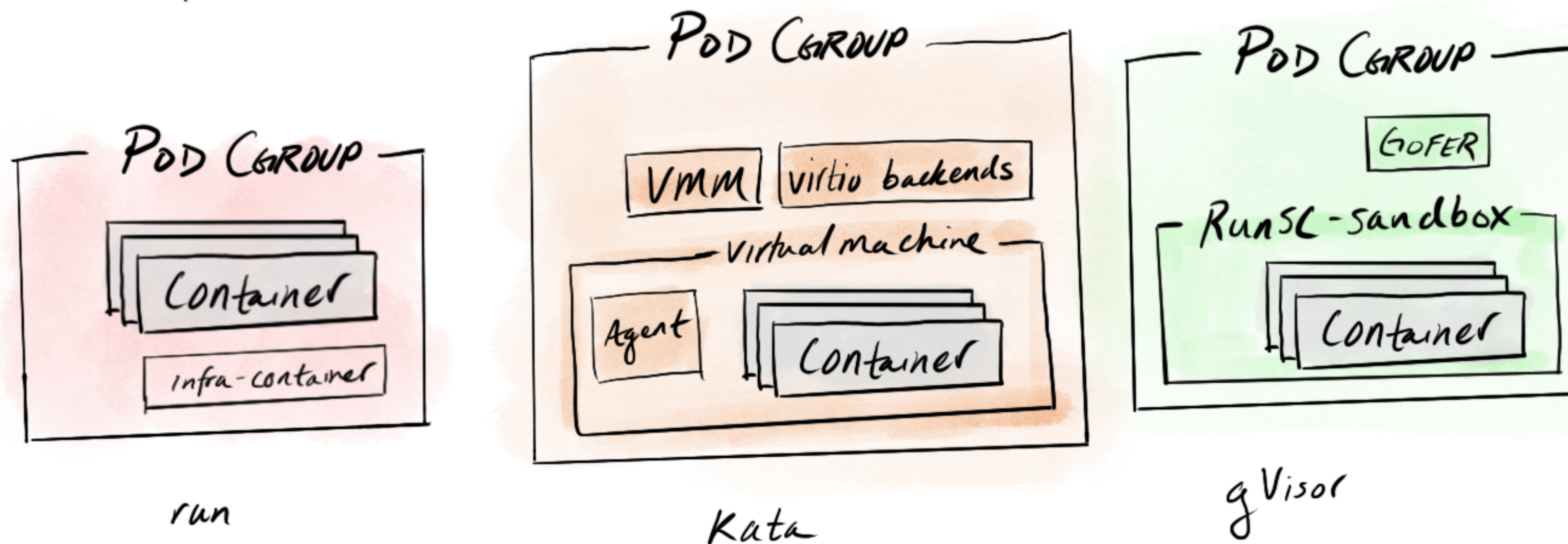


gVisor



Pod Overhead?

Pod Overhead?



pod.resources

!=

sum(container[].resources)

Node Resources

SYSTEM RESERVED

(256 MB)

system services, etc.

KUBE RESERVED

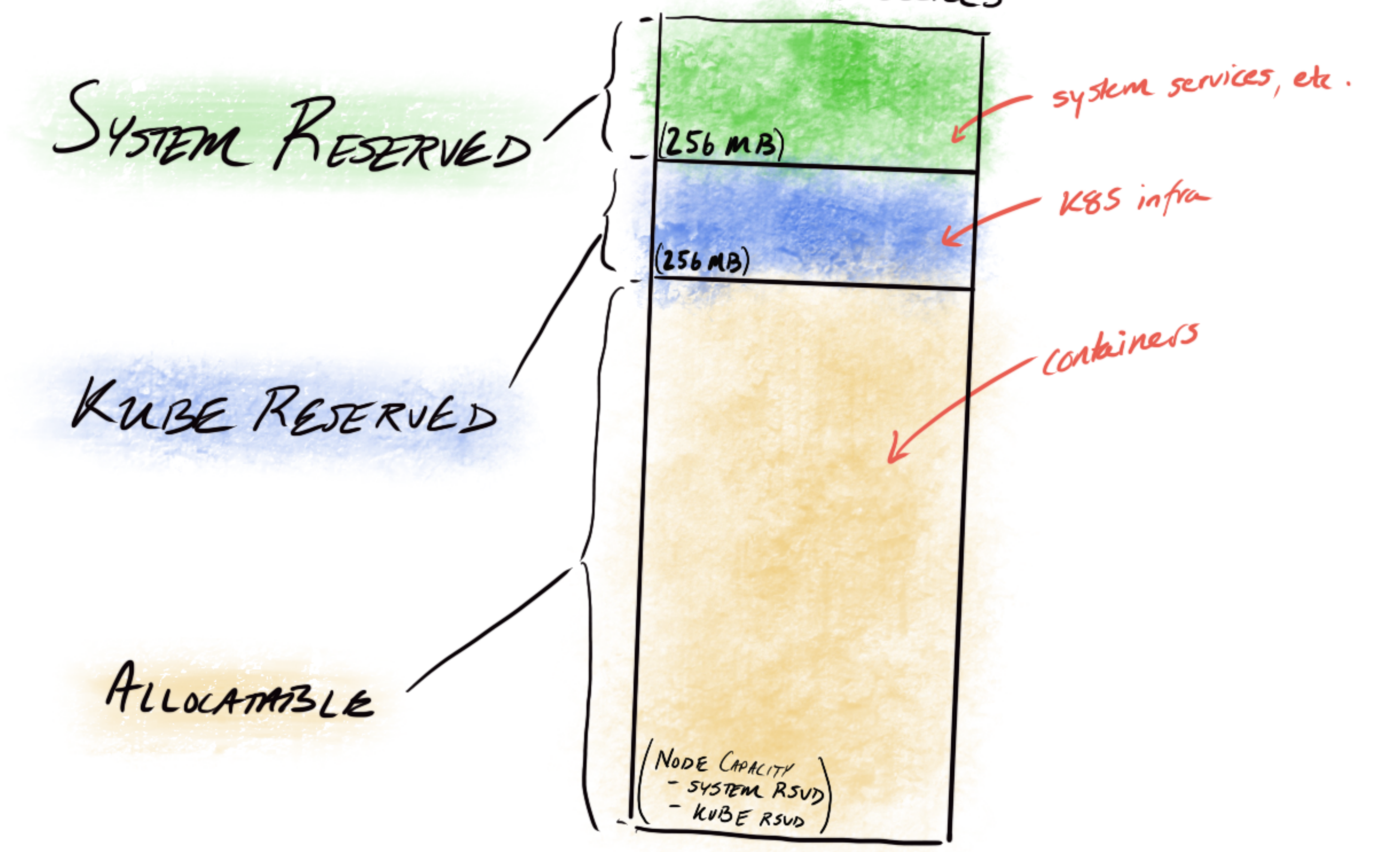
(256 MB)

K8S infra

ALLOCATABLE

containers

(NODE CAPACITY
- SYSTEM RSVD
- KUBE RSVD)



PodOverhead

Alpha feature, introduced in 1.16, which accounts for the overheads associated with running a pod.

How it works

1. Overhead fields are added to RuntimeClass and PodSpec definitions
2. At admission time, update PodSpec to include overhead **iff** a valid overhead is specified in the specified RuntimeClass
3. Account for this overhead in remaining pod life-cycle / management Kubernetes

kubectl get pod busybox-kata-qemu -o yaml

```
1. vim
--
kind: RuntimeClass
apiVersion: node.k8s.io/v1beta1
metadata:
  name: kata-fc
handler: kata-fc
overhead:
  podFixed:
    memory: "130Mi"
    cpu: "250m"
```

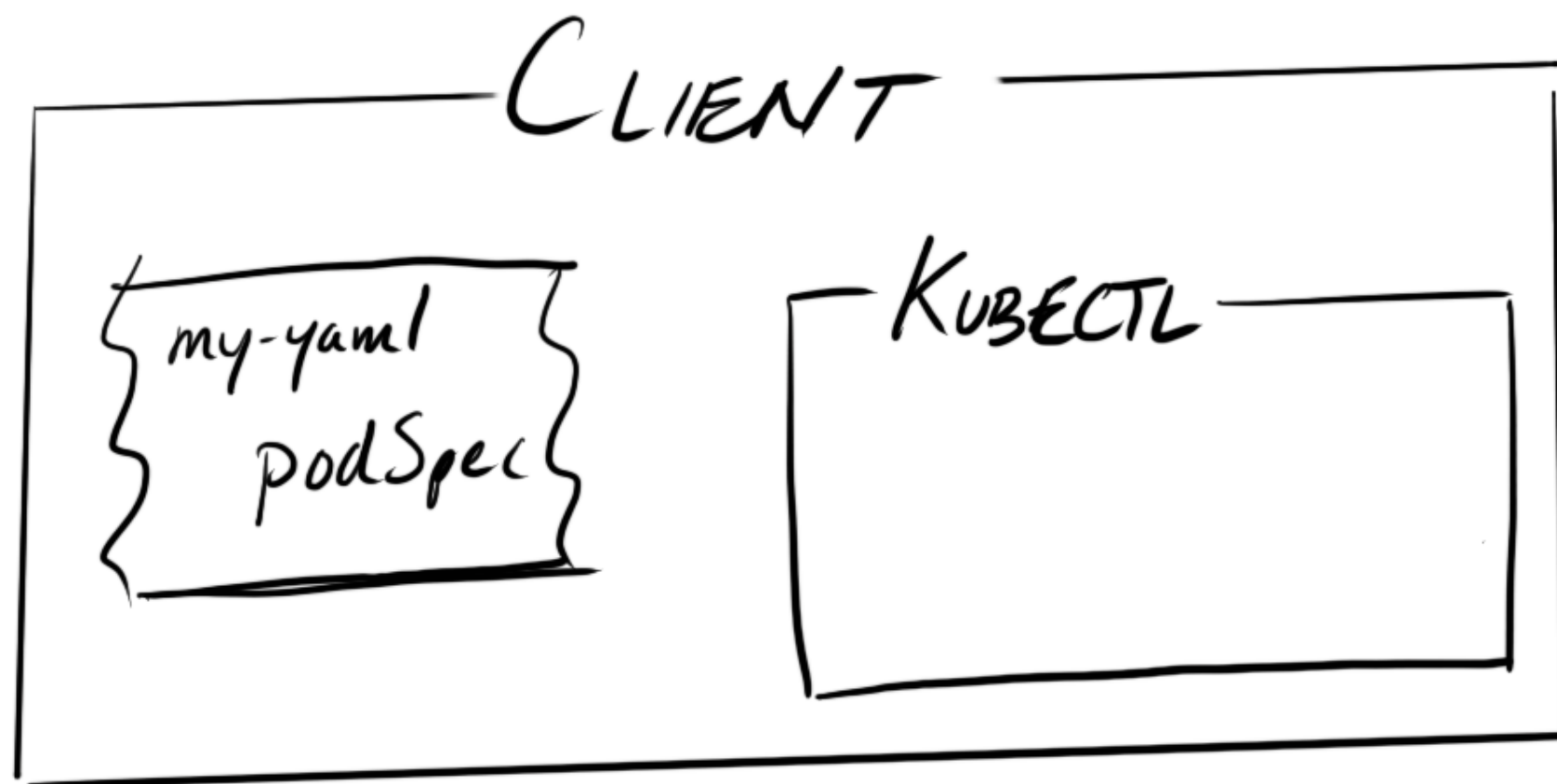
```
1. vim
apiVersion: v1
kind: Pod
metadata: busybox-kata-fc
spec:
  runtimeClassName: kata-fc
  containers:
  - name: busybox-ctr
    image: busybox
    resources:
      limits:
        cpu: 100m
        memory: 100Mi
~
~
```

```
1. k8s@k8s-po: ~ (nc)
limits:
  cpu: 100m
  memory: 100Mi
requests:
  cpu: 100m
  memory: 100Mi
terminationMessagePath: /dev/termination-log
terminationMessagePolicy: File
volumeMounts:
- mountPath: /var/run/secrets/kubernetes.io/serviceaccount
  name: default-token-vw9mh
  readOnly: true
dnsPolicy: ClusterFirst
enableServiceLinks: true
nodeName: k8s-po
overhead:
  cpu: 250m
  memory: 130Mi
priority: 0
restartPolicy: Always
runtimeClassName: kata-fc
schedulerName: default-scheduler
securityContext: {}
serviceAccount: default
serviceAccountName: default
terminationGracePeriodSeconds: 30
tolerations:
-- VISUAL --
[0] 0:vi*Z "k8s-po" 04:26 18-Nov-19
```

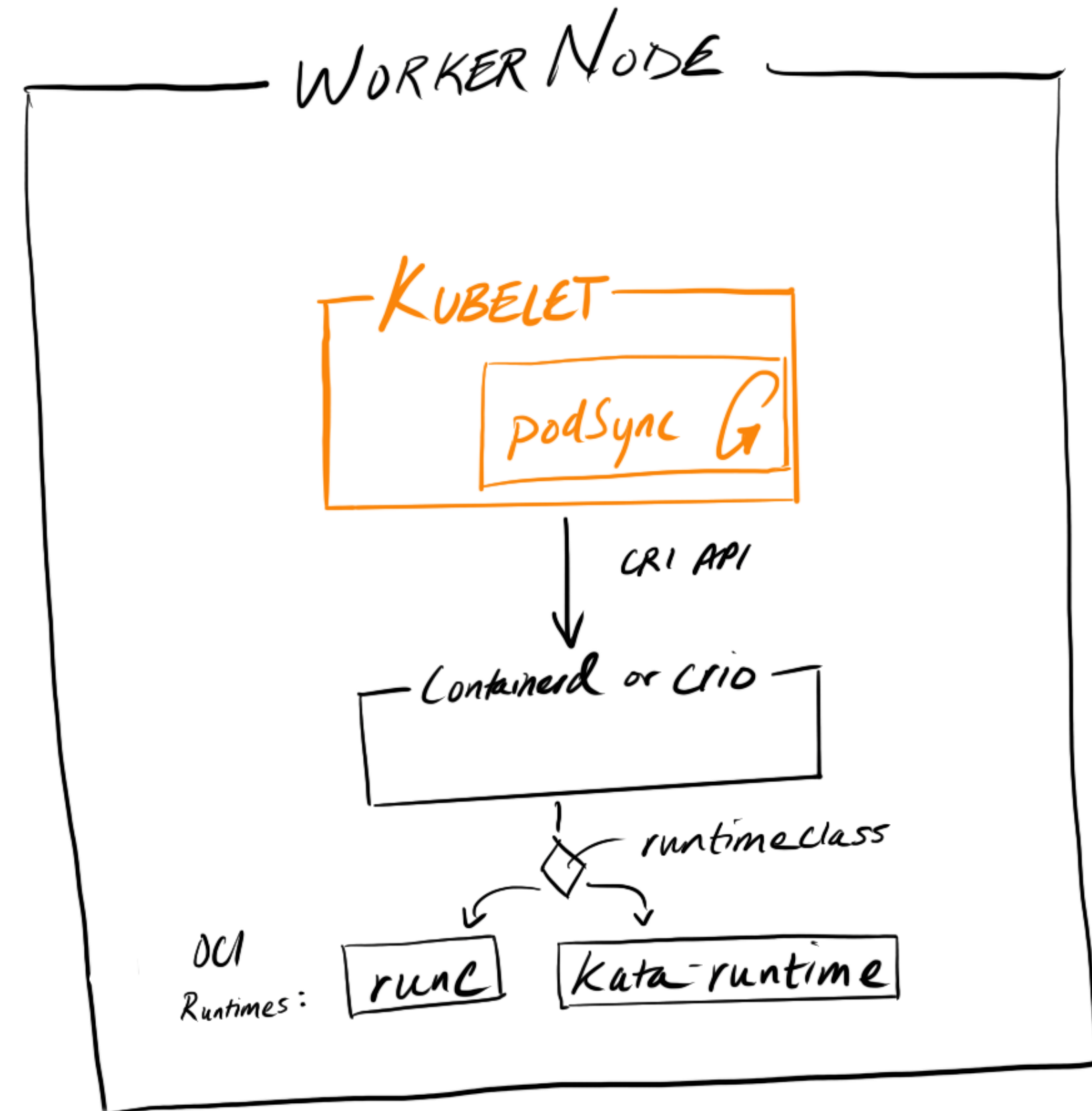
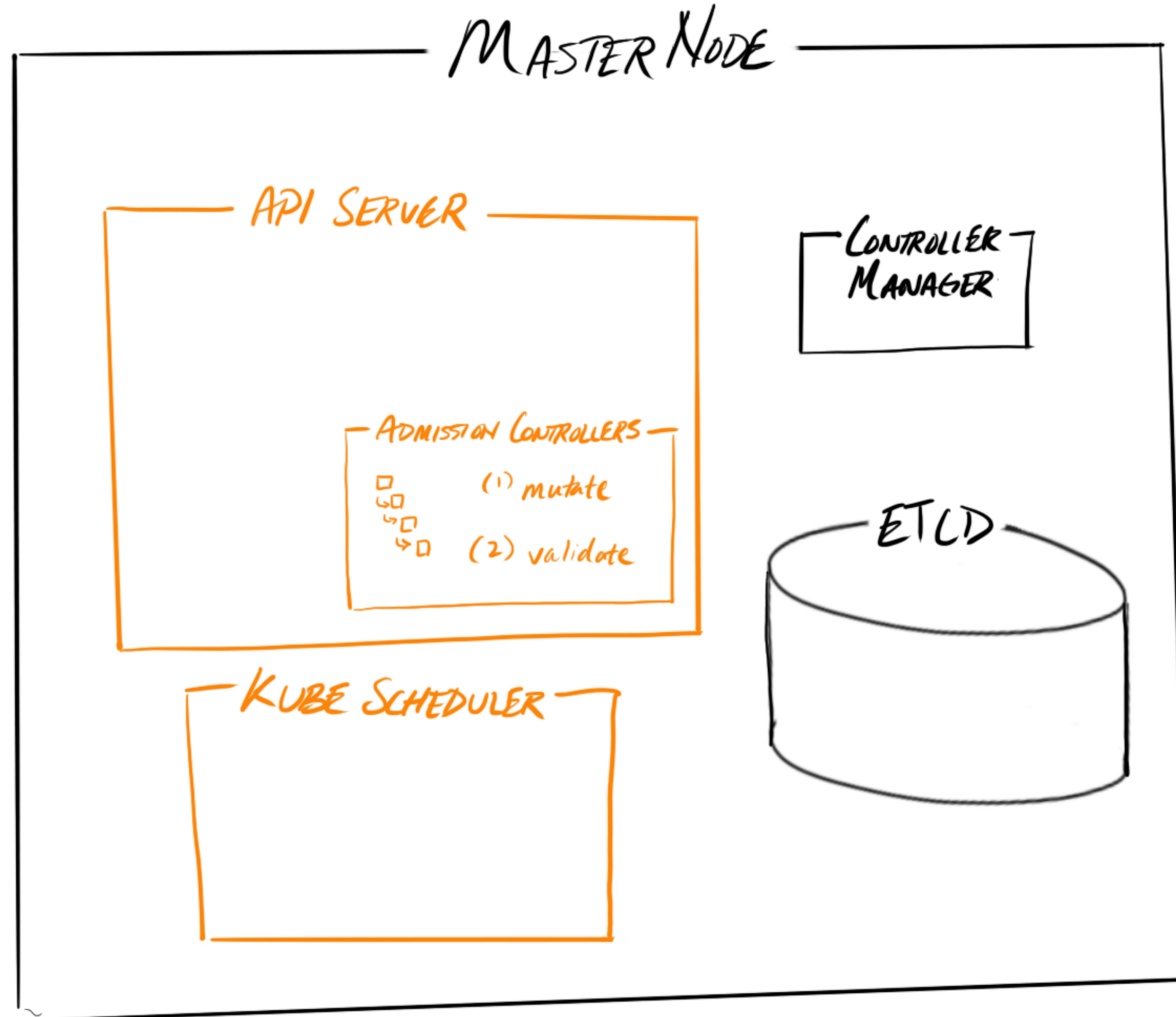
```
1. vim 🛎️
--
kind: RuntimeClass
apiVersion: node.k8s.io/v1beta1
metadata:
  name: kata-fc
handler: kata-fc
overhead:
  podFixed:
    memory: "130Mi"
    cpu: "250m"
```

```
1. vim
apiVersion: v1
kind: Pod
metadata: busybox-kata-fc
spec:
  runtimeClassName: kata-fc
  containers:
  - name: busybox-ctr
    image: busybox
    resources:
      limits:
        cpu: 100m
        memory: 100Mi
~
~
```

```
1. k8s@k8s-po: ~ (nc)
k8s@k8s-po: $ cat /sys/fs/cgroup/memory/kubepods/pod*/memory.limit_in_bytes
241172480
k8s@k8s-po: $ cat /sys/fs/cgroup/cpu,cpuacct/kubepods/pod*/cpu.cfs_quota_us
35000
5,0-1 Top
[0] 0:vi*Z "k8s-po" 04:40 18-Nov-19
```

PodSpec API (core)
 RuntimeClass API (node)



Changes Required

```
1. ernst@ernstworkstation: ~/go/src/k8s.io/kubernetes (ssh)
...amples (bash) 1  ...tation: ~ (bash) 2  ...tation: ~ (bash) 3  vim 4  ...o/kubernetes (ssh) 5  bash 6
56
57 // PodRequestsAndLimits returns a dictionary of all defined resources summed up for all
58 // containers of the pod. If PodOverhead feature is enabled, pod overhead is added to the
59 // total container resource requests.
60 func PodRequestsAndLimits(pod *v1.Pod) (reqs, limits v1.ResourceList) {
61     reqs, limits = v1.ResourceList{}, v1.ResourceList{}
62     for _, container := range pod.Spec.Containers {
63         addResourceList(reqs, container.Resources.Requests)
64         addResourceList(limits, container.Resources.Limits)
65     }
66     // init containers define the minimum of any resource
67     for _, container := range pod.Spec.InitContainers {
68         maxResourceList(reqs, container.Resources.Requests)
69         maxResourceList(limits, container.Resources.Limits)
70     }
71
72     // if PodOverhead feature is supported, add overhead for running a pod
73     // to the sum of container requests:
74     if pod.Spec.Overhead != nil && utilfeature.DefaultFeatureGate.Enabled(features.PodOverhead) {
75         addResourceList(reqs, pod.Spec.Overhead)
76     }
77
78     return
79 }
80
```

80,0-1 27%

Who cares?

Users

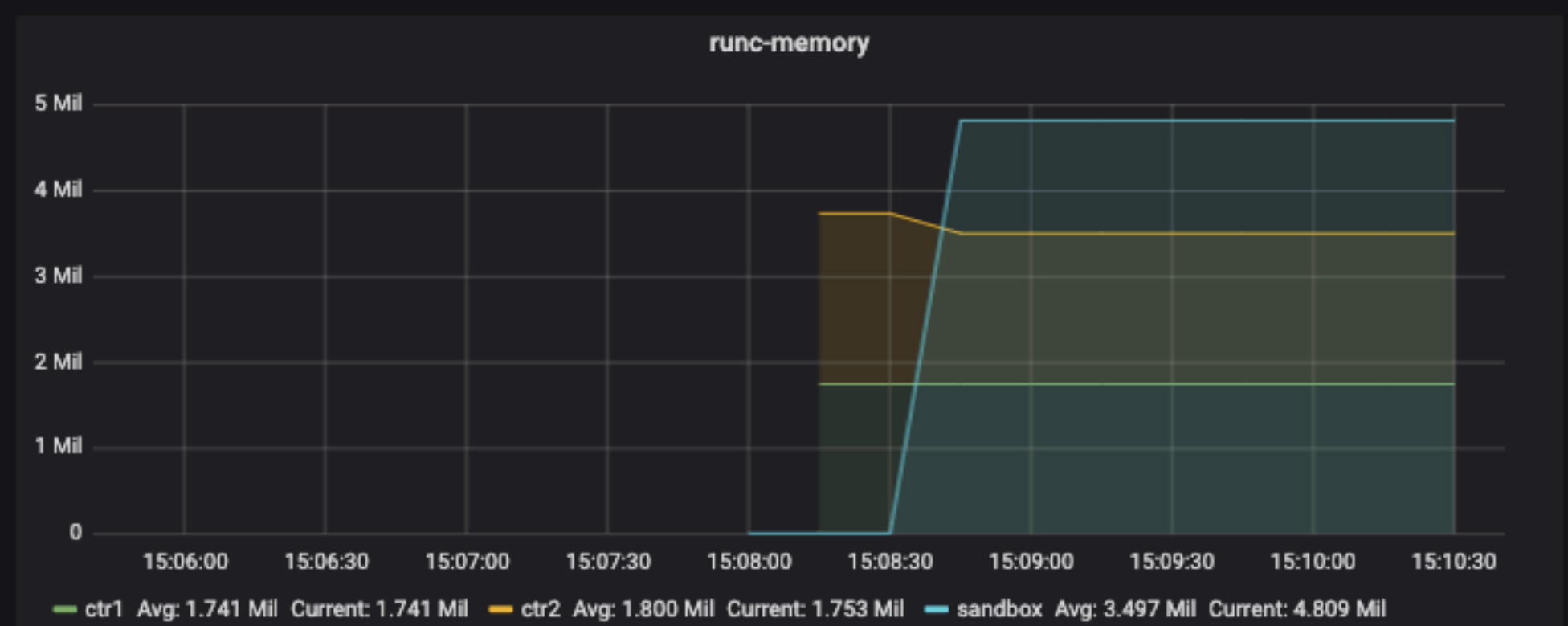
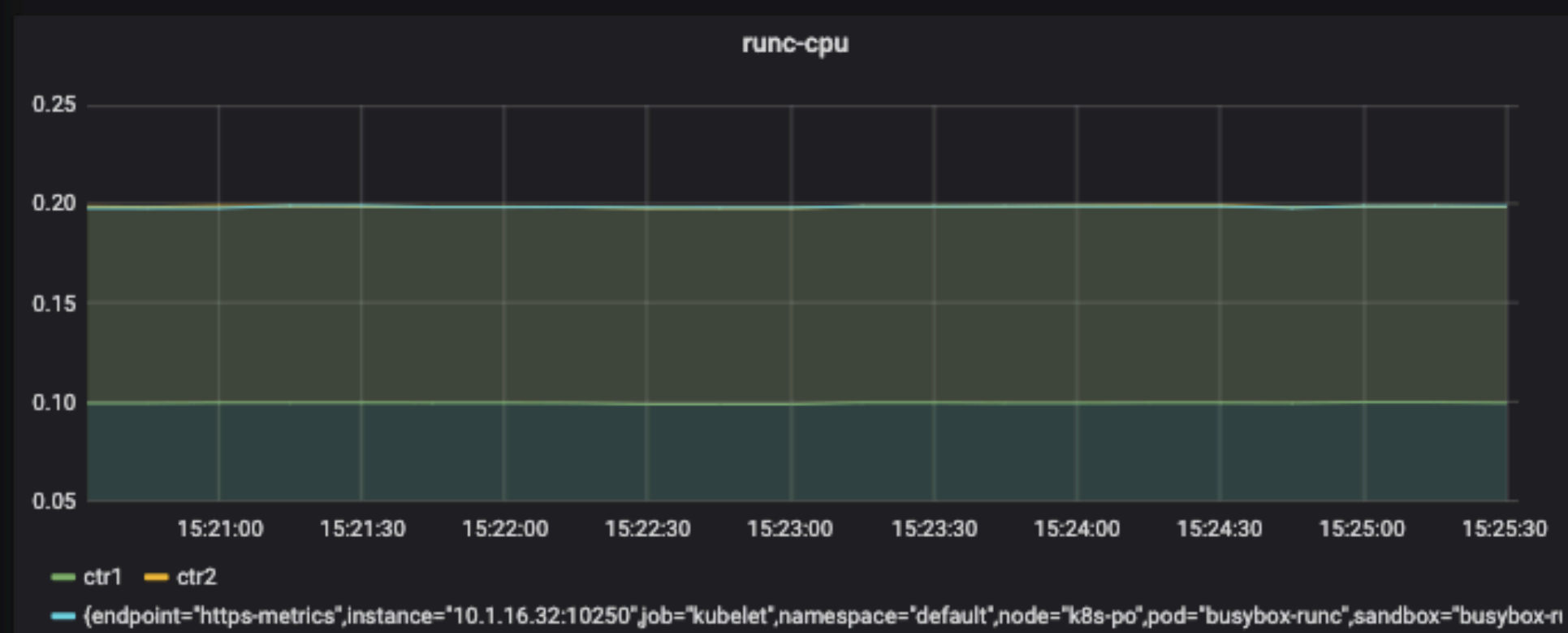
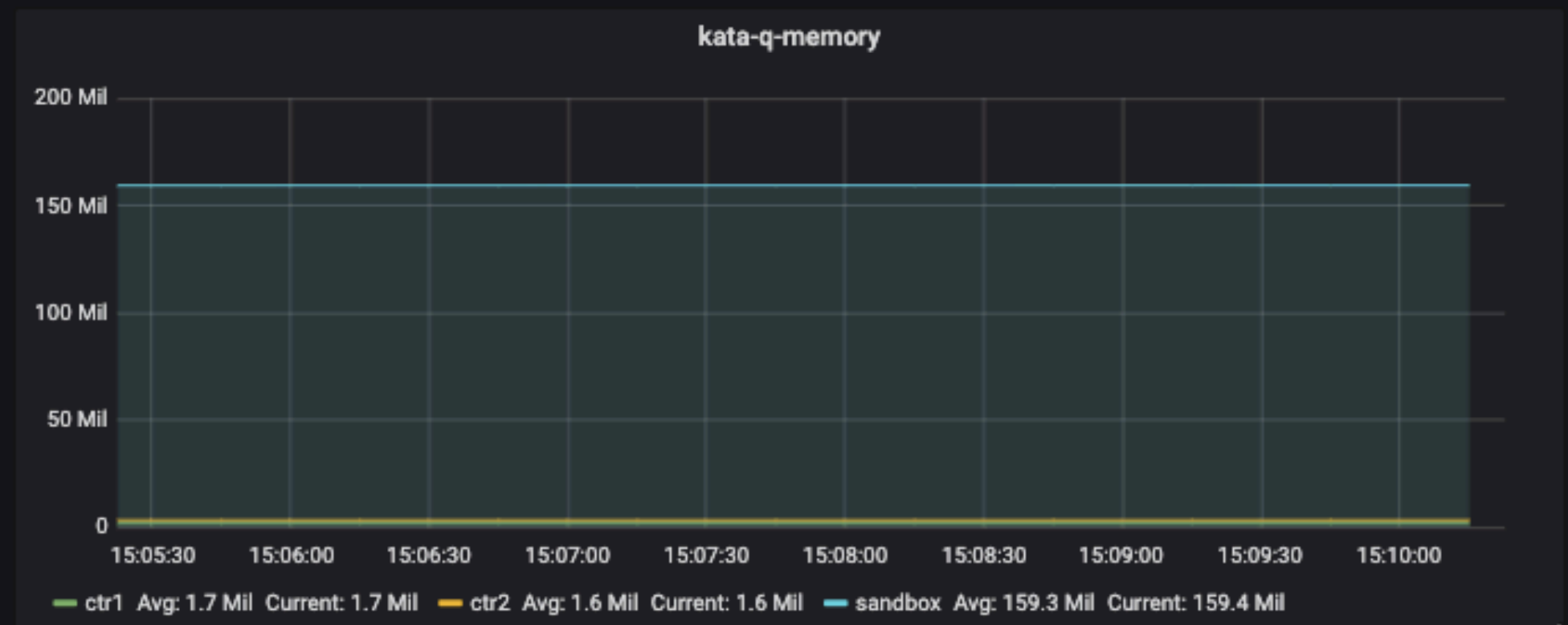
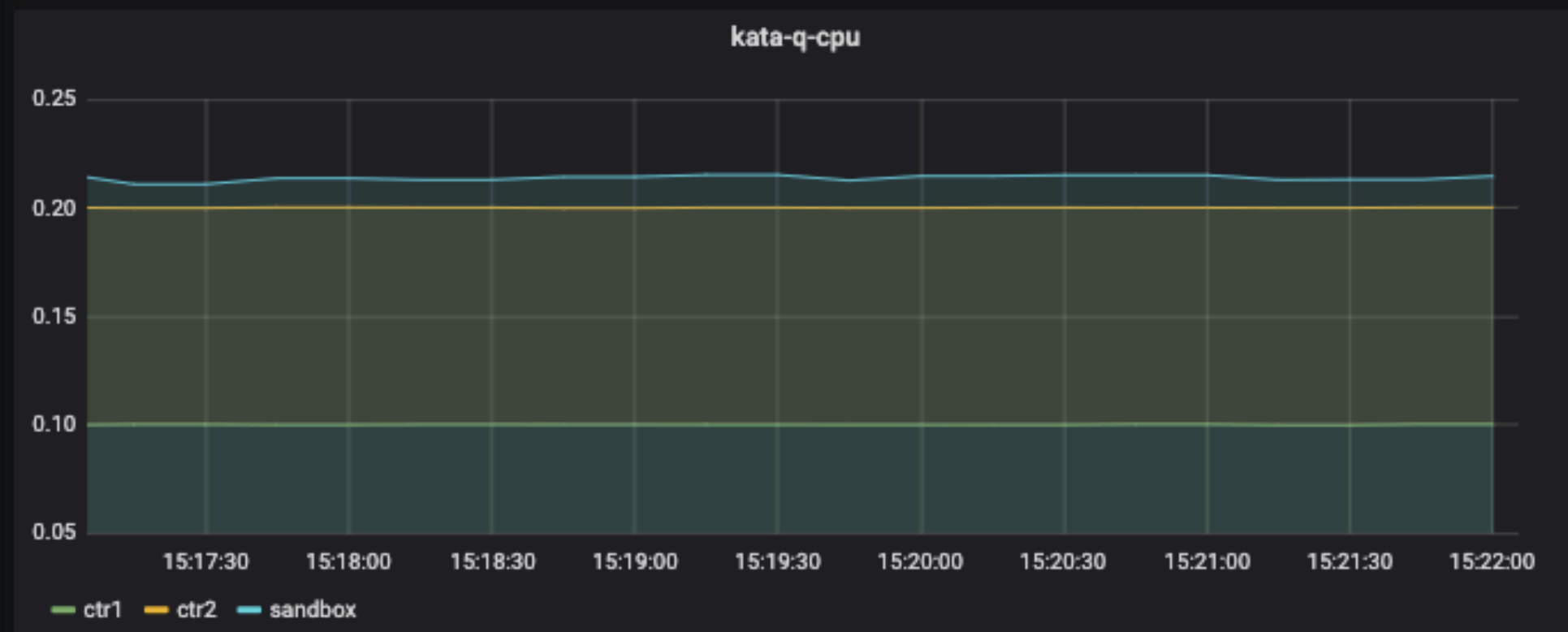
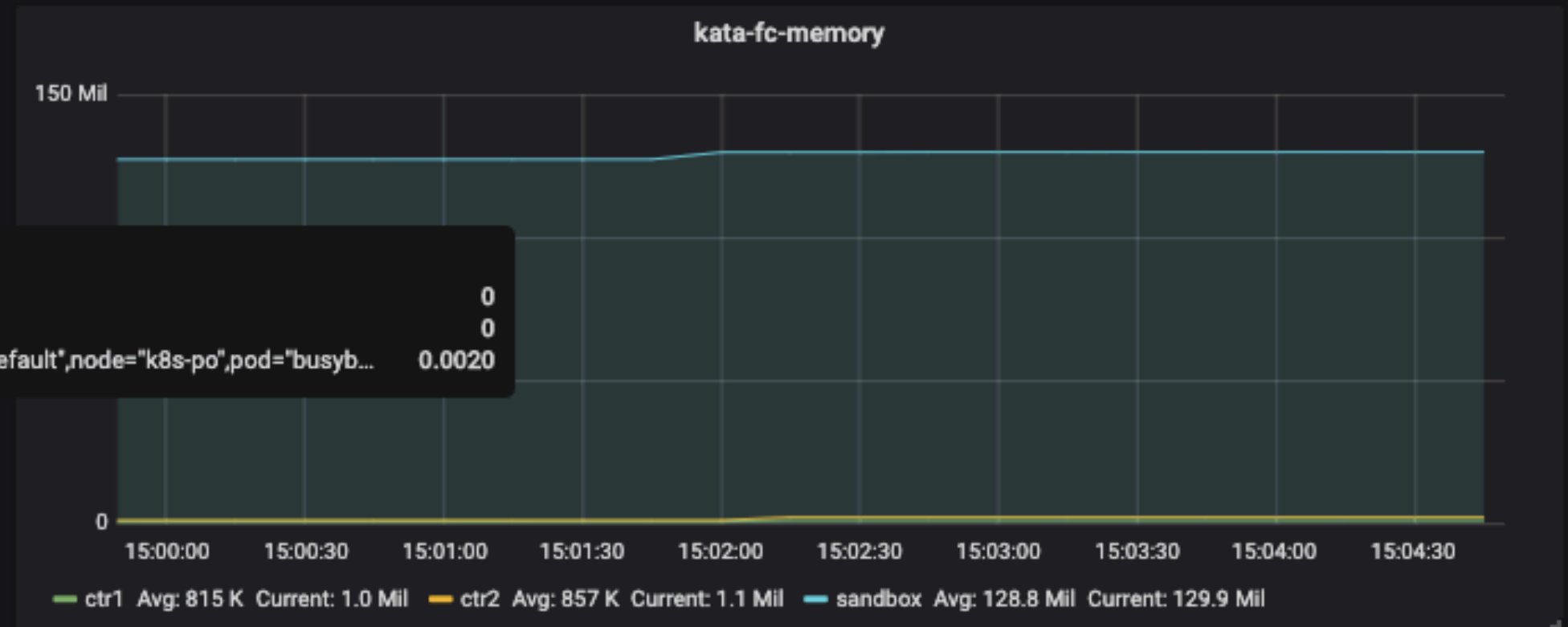
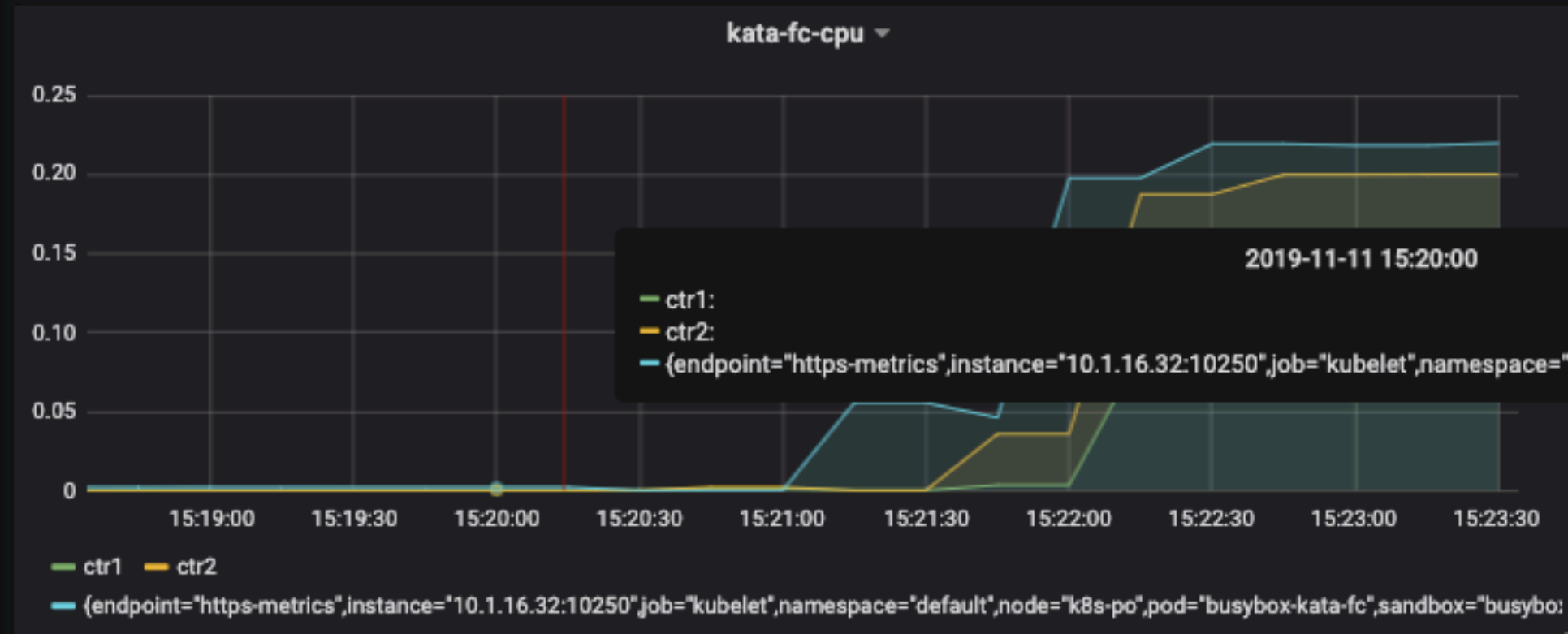
Severely over-constrained: Oom...Why isn't my workload running?

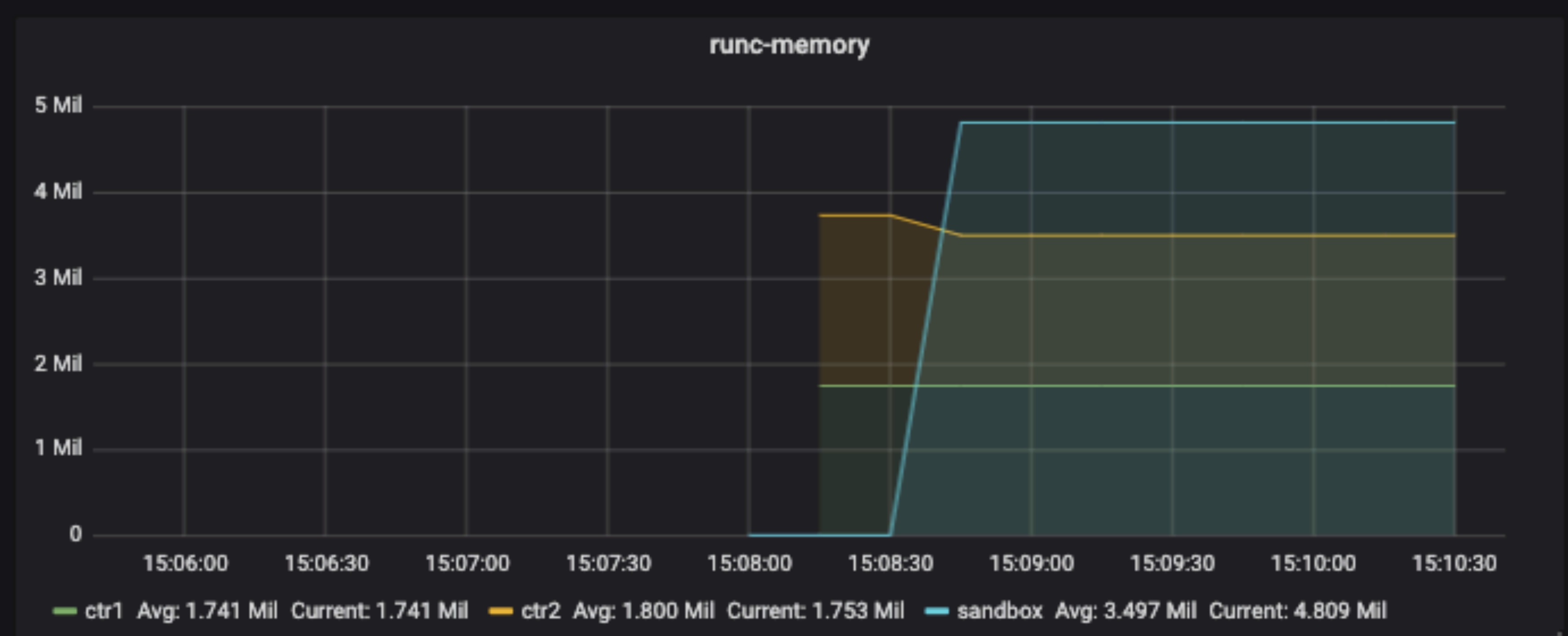
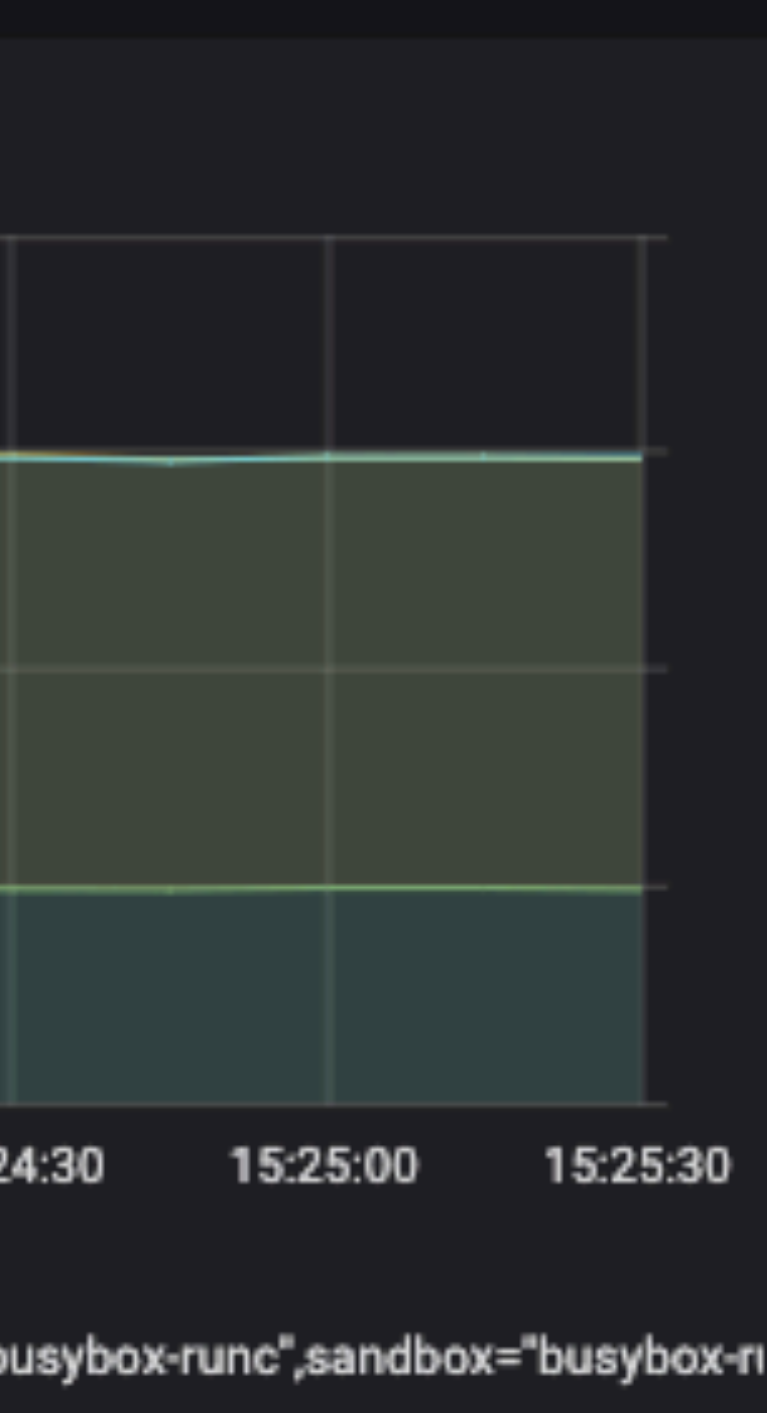
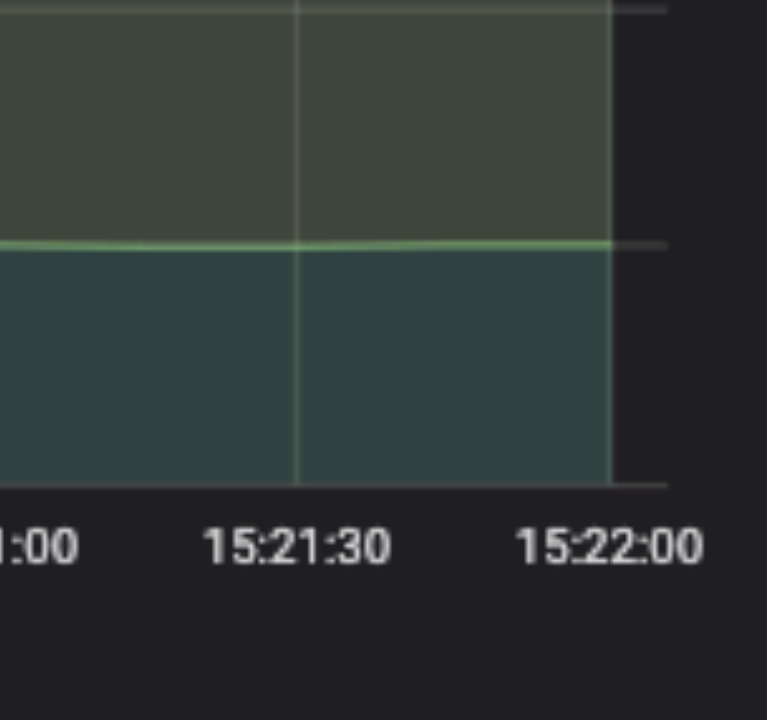
over-constrained: inconsistent, poor performance

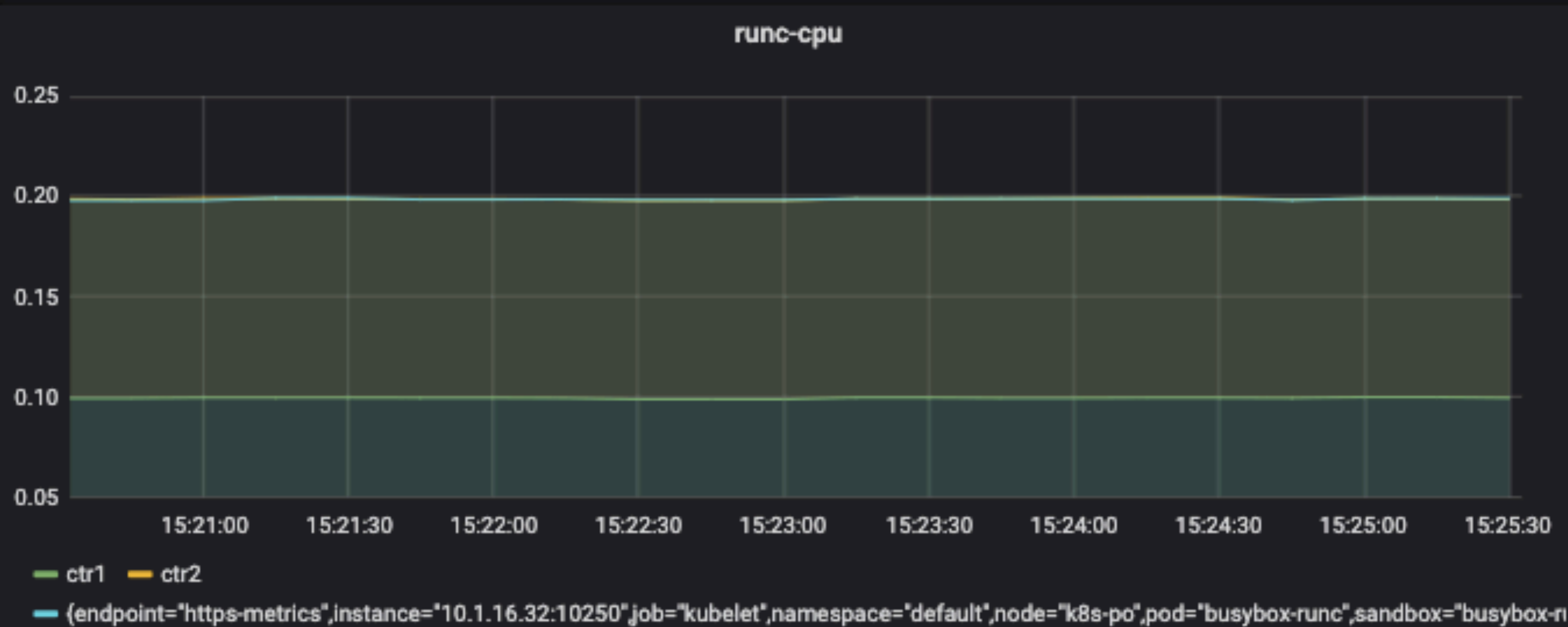
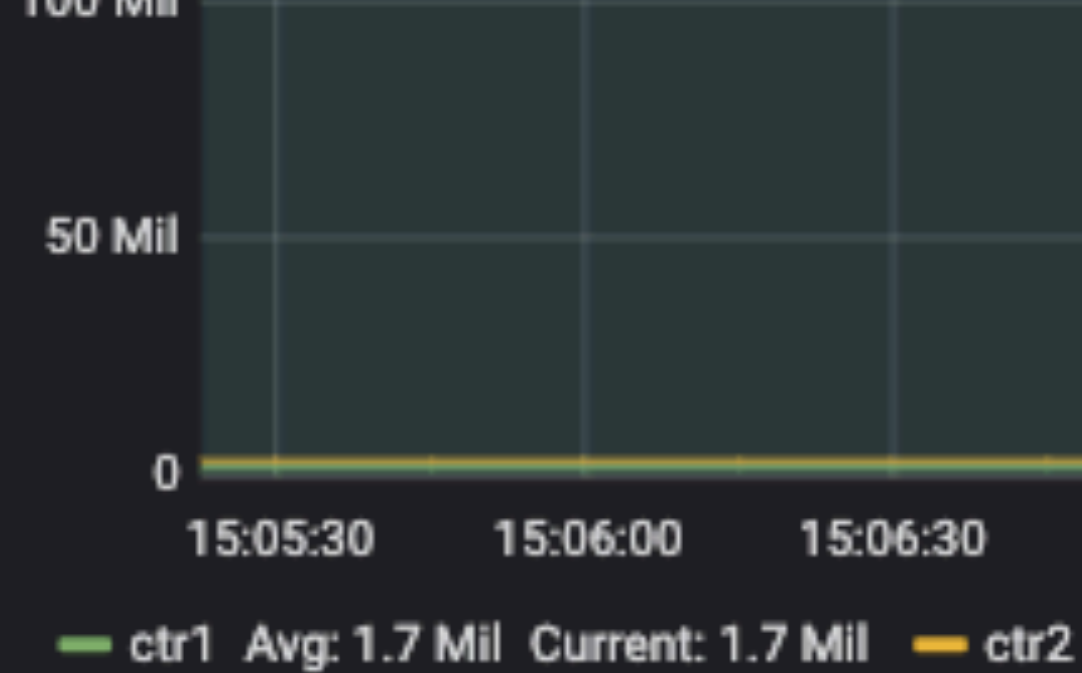
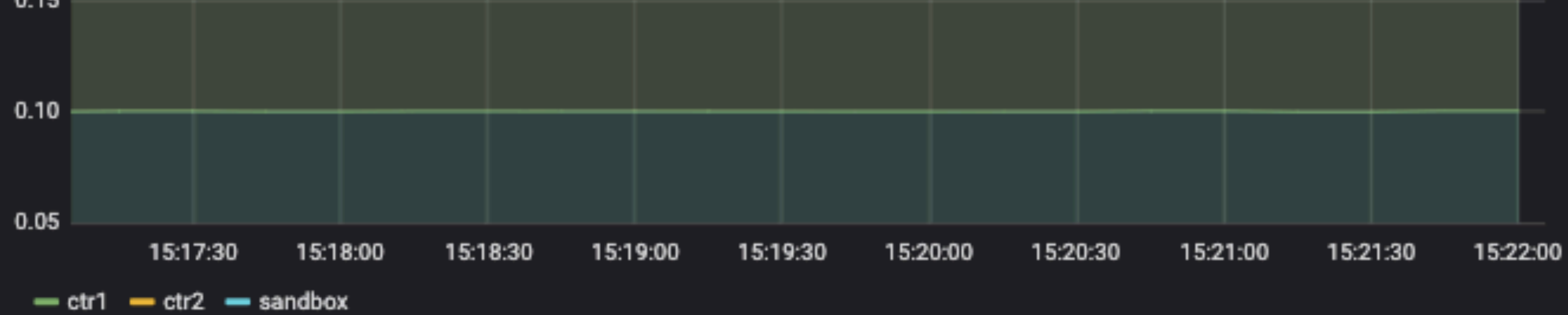
```
k8s@k8s-po:~/k8s-pod-overhead/test-workloads$  
k8s@k8s-po:~/k8s-pod-overhead/test-workloads$  
k8s@k8s-po:~/k8s-pod-overhead/test-workloads$ █
```

Administrators

ResourceQuota doesn't reflect reality

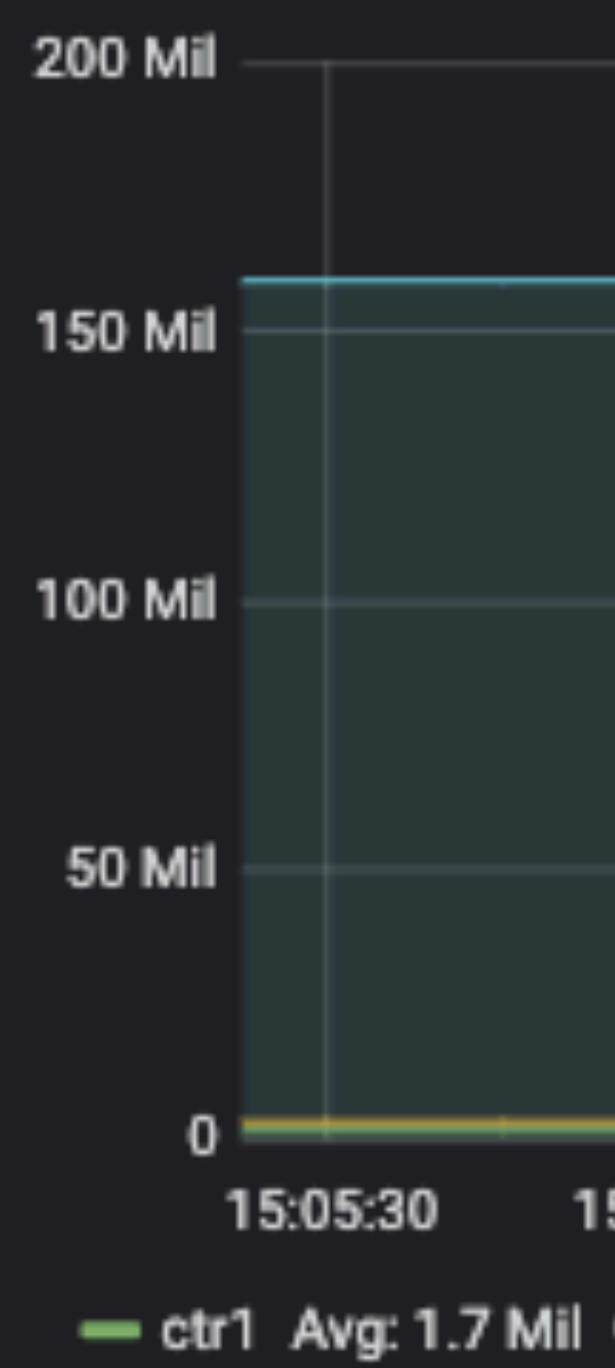
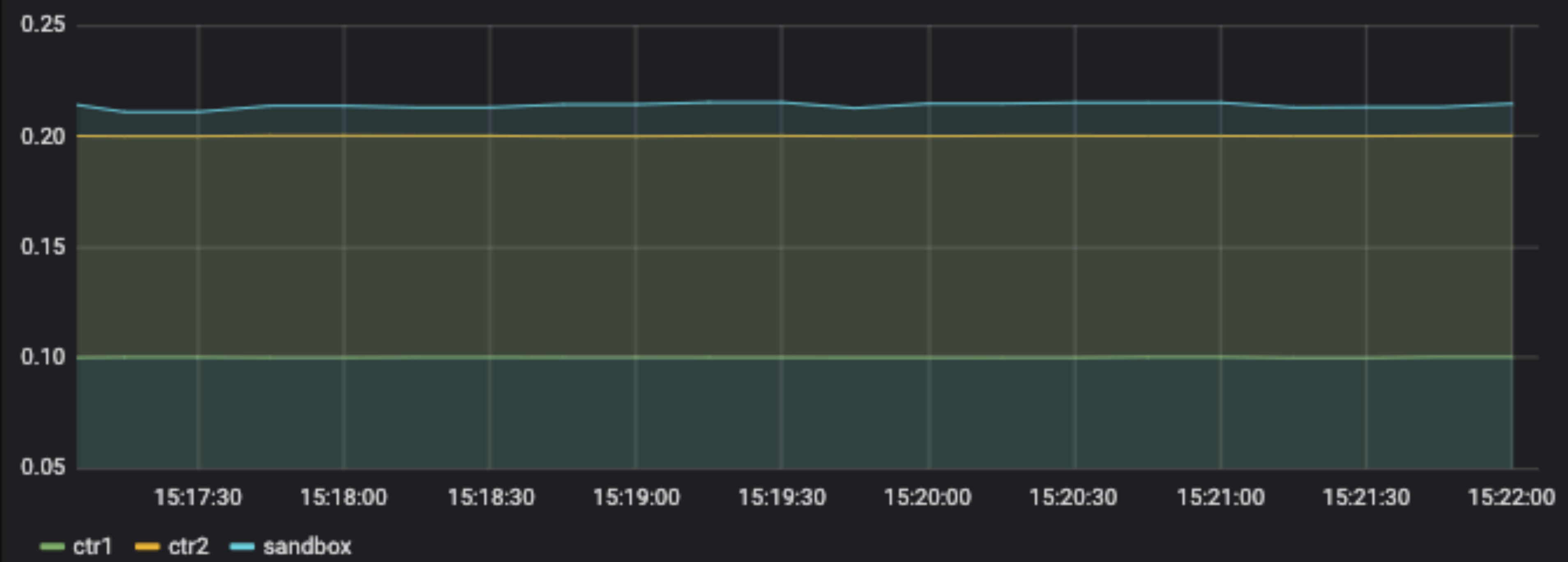




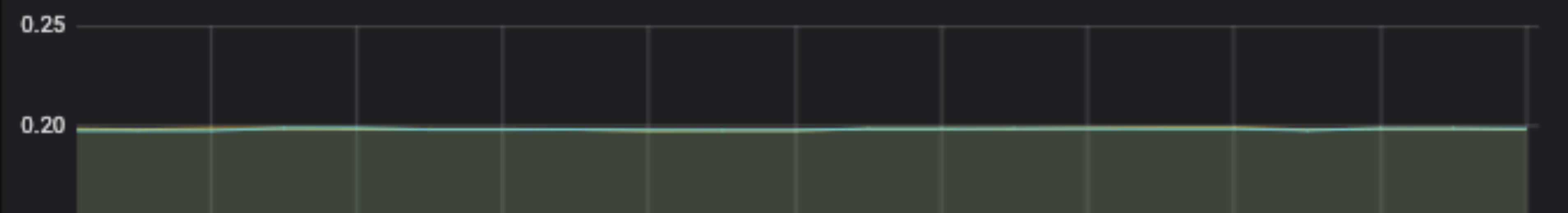


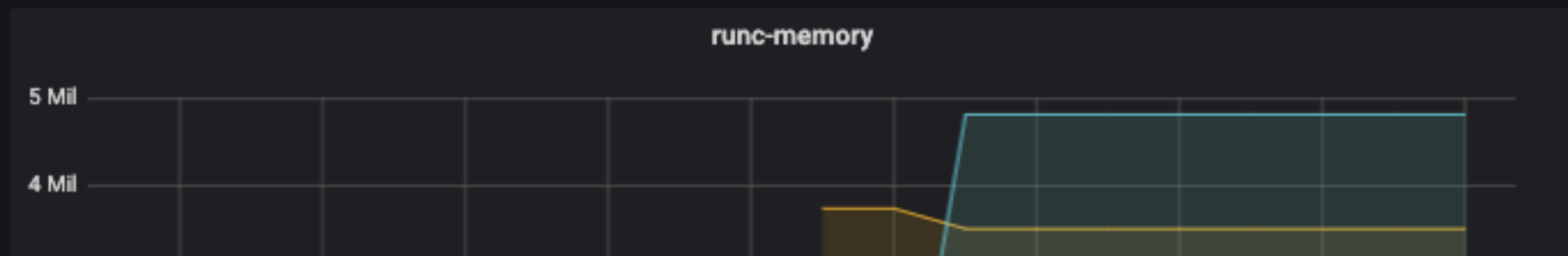
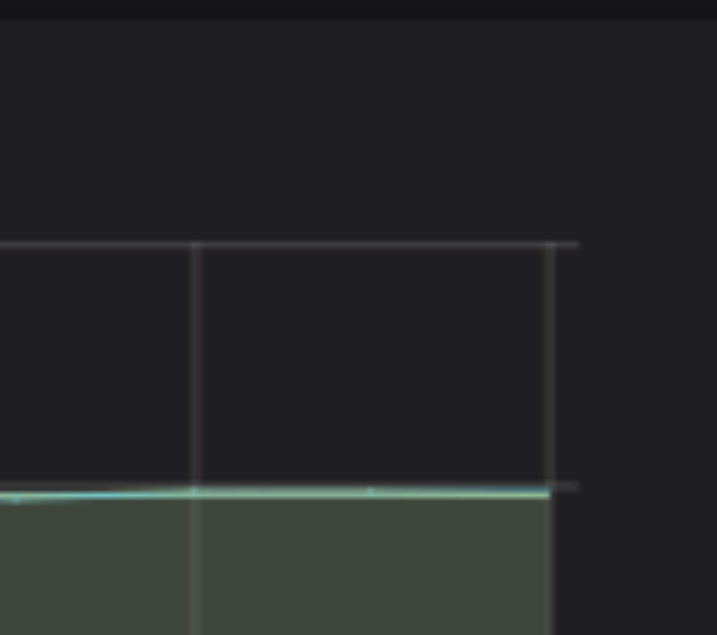
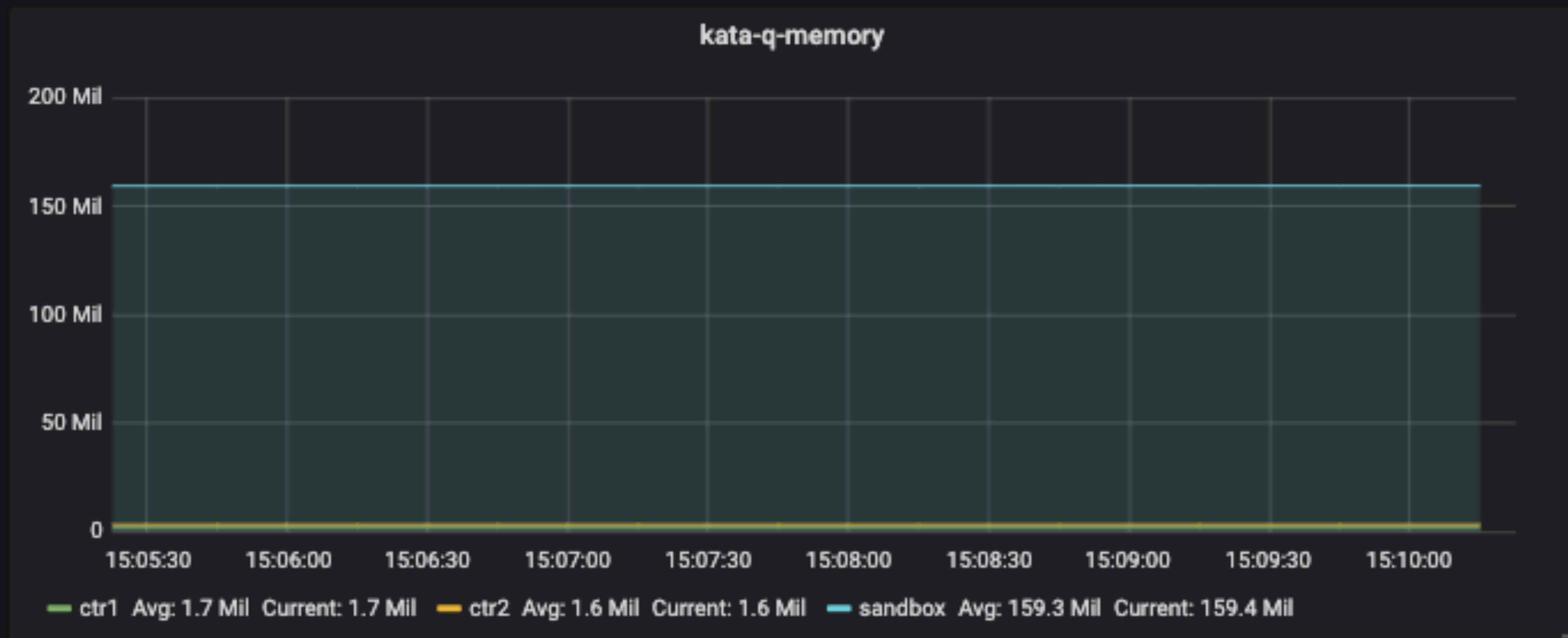
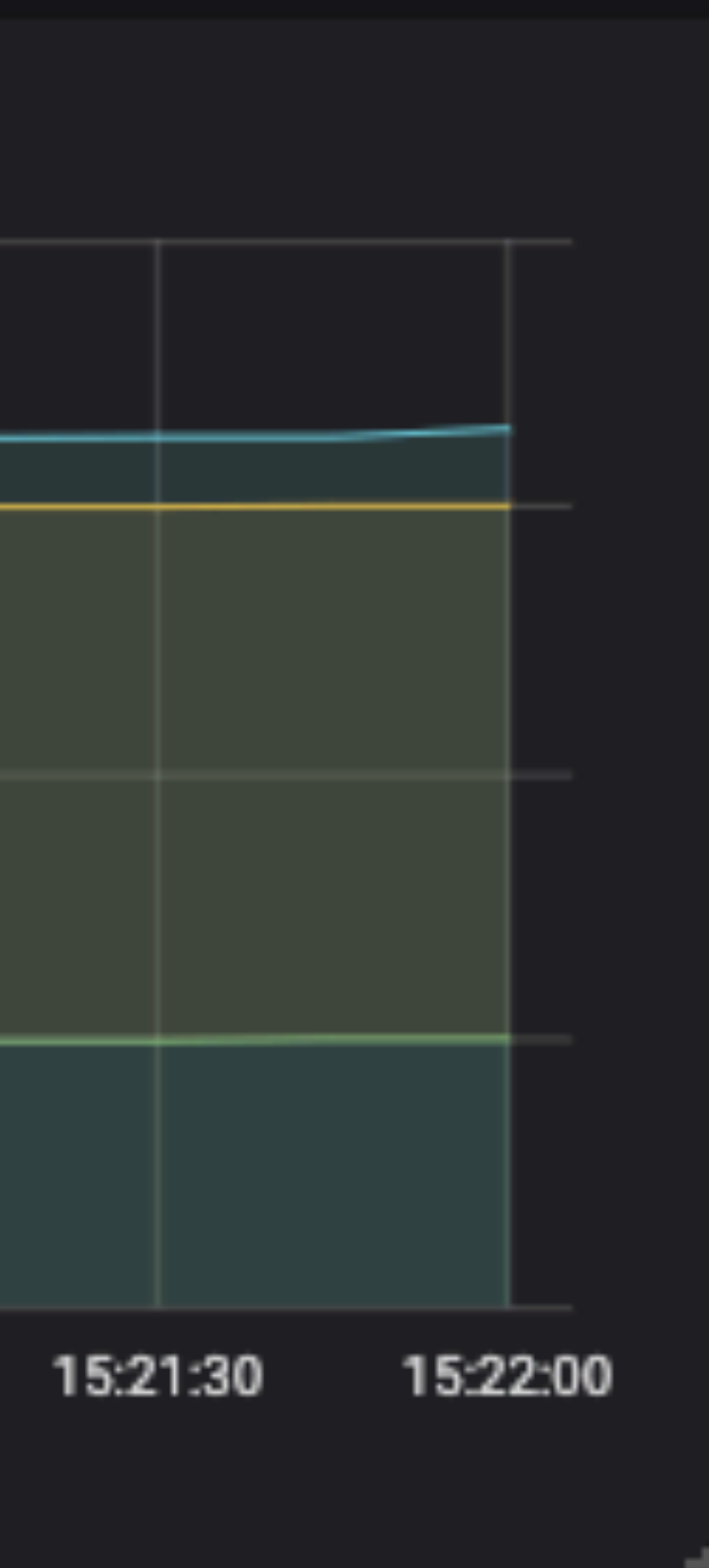
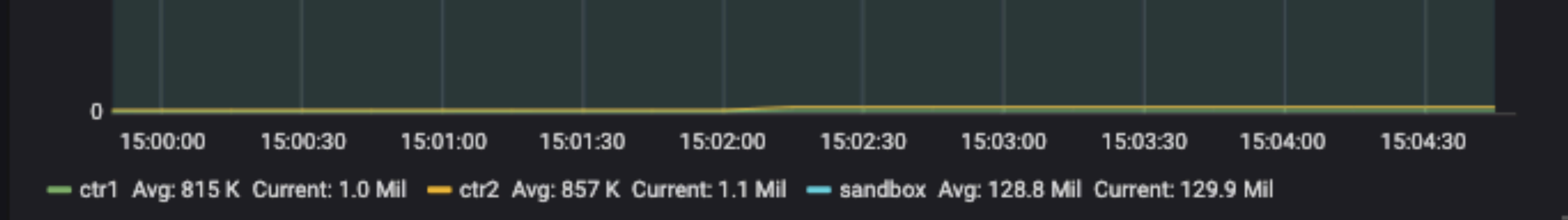
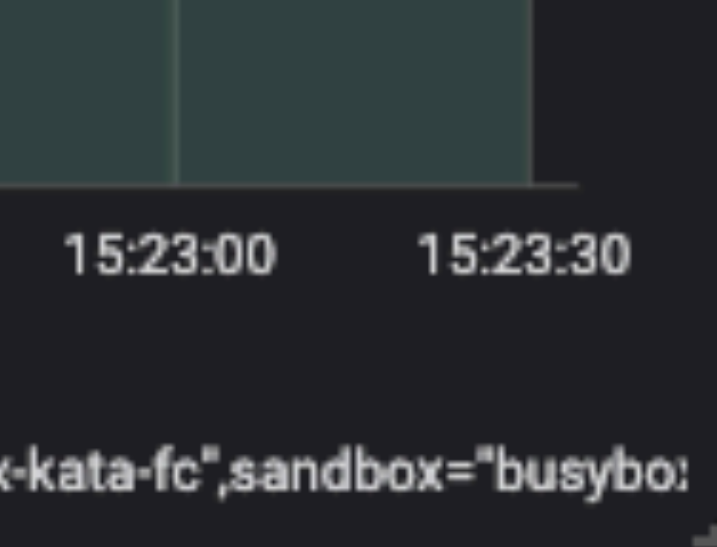


kata-q-cpu

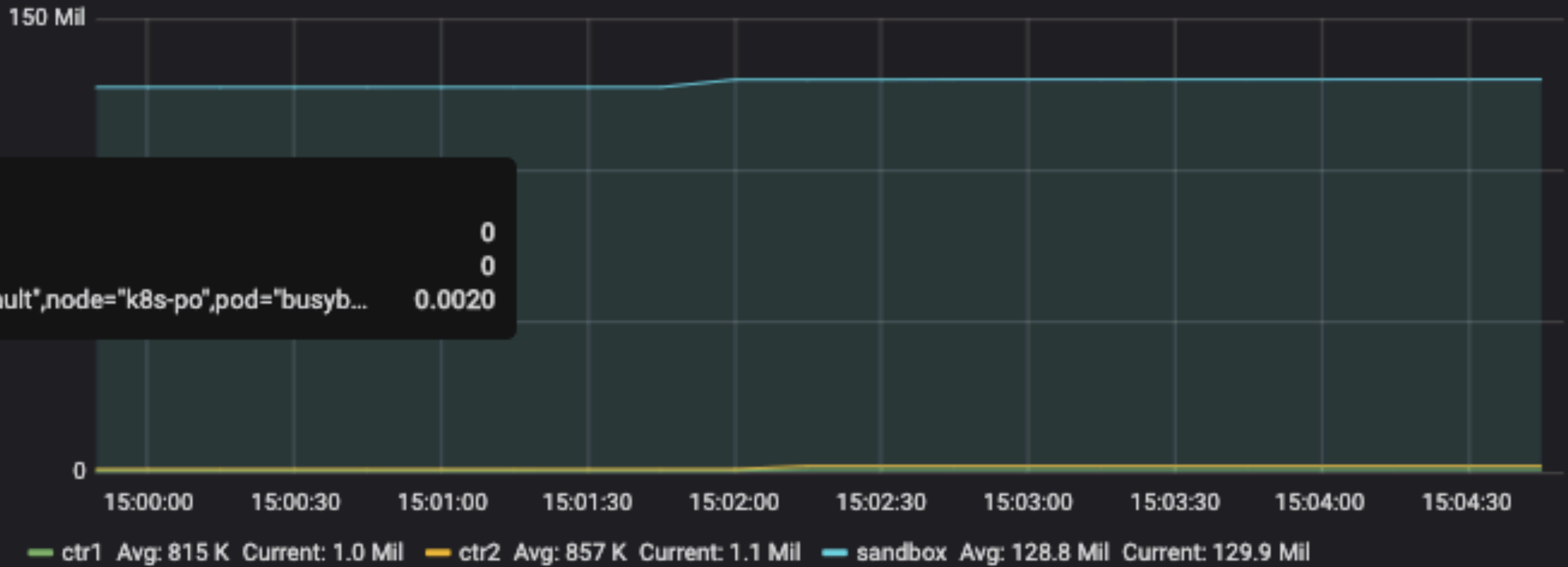


runc-cpu

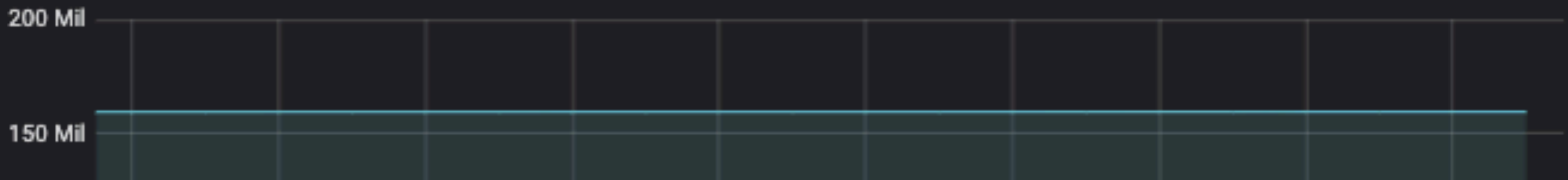




kata-fc-memory



kata-q-memory



Providers

No room in the pod slice? We'll just drop xyz into system slice, and hope that's sized okay

If we don't account, who pays for it?

PodOverhead status

- Available in Kubernetes 1.16 as an alpha feature
- Overheads are static - good starting point, but not always realistic
- Expected to move to beta in 1.18 release

```
eernst — k8s@k8s-po: ~/k8s-pod-overhead — nc ◀ ssh k8s@40.65.124.189 — 61x33
apiVersion: kubeadm.k8s.io/v1beta1
kind: InitConfiguration
nodeRegistration:
  criSocket: "/var/run/containerd/containerd.sock"
  kubeletExtraArgs:
    feature-gates: PodOverhead=true
---
apiVersion: kubelet.config.k8s.io/v1beta1
kind: KubeletConfiguration
cgroupDriver: cgroupfs
featureGates:
  PodOverhead: true
systemReserved:
  cpu: 500m
  memory: 256M
kubeReserved:
  cpu: 500m
  memory: 256M
---
apiVersion: kubeadm.k8s.io/v1beta1
kind: ClusterConfiguration
networking:
  dnsDomain: cluster.local
  podSubnet: 10.244.0.0/16
  serviceSubnet: 10.96.0.0/12
apiServer:
  extraArgs:
    feature-gates: PodOverhead=true
scheduler:
  extraArgs:
    feature-gates: PodOverhead=true
~
6,20 All
```





KubeCon



CloudNativeCon

North America 2019

Thanks!



Learnings

- API Reviews can be painful, and are the longest pole in the tent
- Kubernetes makes use of a lot of auto-generated machinery, especially for API version conversions
- Feature wise, it is easy to get changes in Kubernetes which will improve node stability and/or more accurate accounting
- Time spent { coding << herding cats for reviews }
- Time spent coding:
 - { writing unit tests and fixing > tracing code to see what happens where >> writing actual feature }