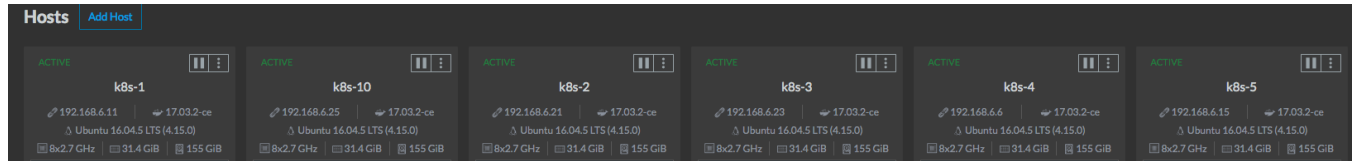


# Kubernetes cluster stability

ONAP has grown in Casablanca and requires a lot of resources to run. Although I have a 6 node K8s cluster that I deployed using rancher, I believe the same kubelet flags can be applied to any method of spinning up and maintaining a K8s cluster (kubeadm etc). The nodes are Openstack VMs that I spun up on a private lab network.



The flavour I used:

VCPUs	RAM	Root Disk
8	32GB	160GB

That gives me 42 cores of CPU, 192 GB RAM, 960 GB storage with which to "try" and deploy all of ONAP.

Out of the box, what I found was that after a few days one of my nodes would become unresponsive and disconnect from the cluster. Even accessing the console through Openstack was slow. The resolution for me was a hard shutdown from the Openstack UI which brought it back to life (until it eventually happened again!).

It seems that the node that died was overloaded with pods that starved the core containers that maintain connectivity to the K8s cluster (kubelet, rancher-agent, etc).

The solution I found was to use kubelet eviction flags to protect the core containers from getting starved of resources.



These flags will not enable you to run all of ONAP on less resources. They are there to simply prevent the overloading of your K8s cluster from crashing. If you attempt to deploy too much the pods will be evicted by K8s.

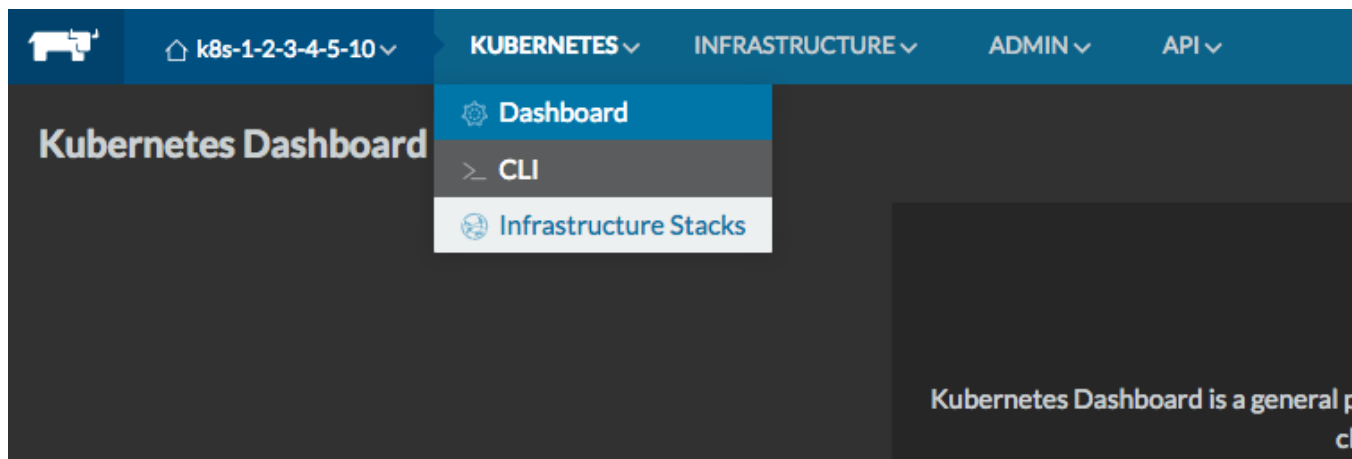
The [resource limit story](#) introduced in Casablanca add minimum requirements to each application. This is a WIP and will help the community better understand how many actual CPU cores and RAM is really required in your K8s cluster to run things.

For more information on the kubelet flags see:

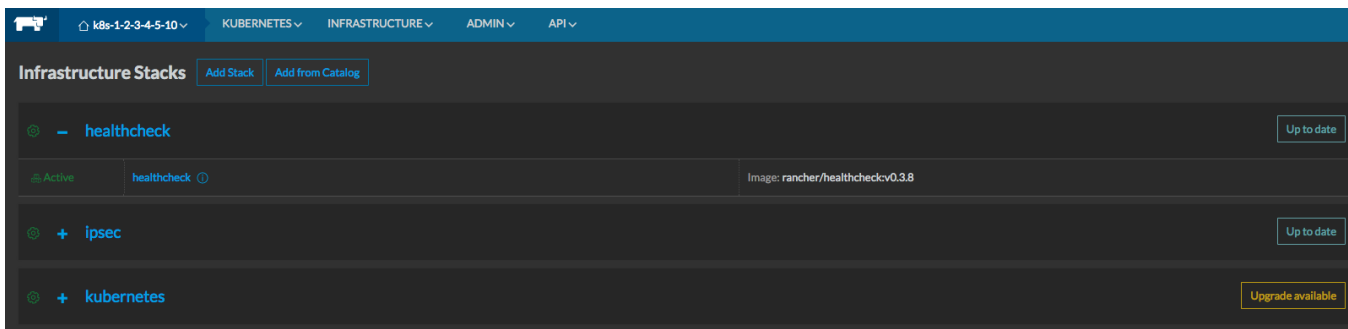
<https://kubernetes.io/docs/tasks/administer-cluster/out-of-resource/>

Here is how I applied them to my already running K8s cluster using rancher:

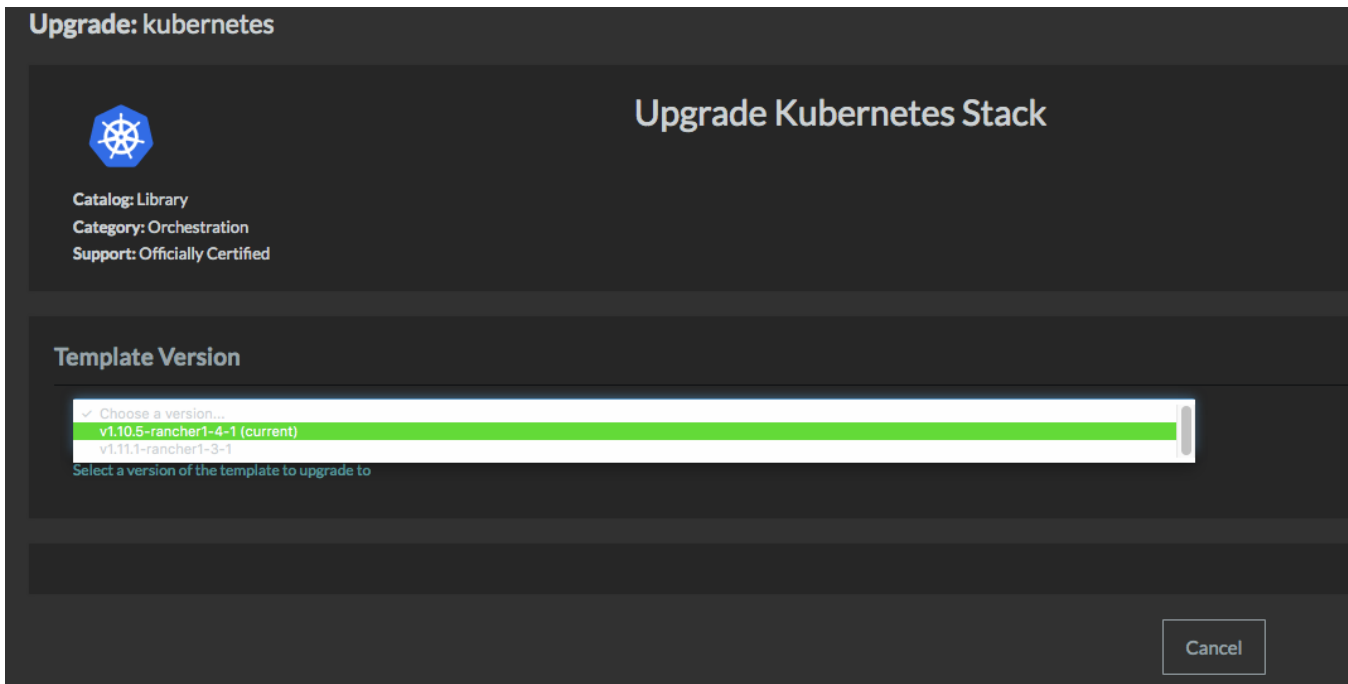
- Navigate to the Kubernetes > "Infrastructure Stacks" menu option



- Depending on the version of rancher and the K8s template version you have deployed, the button on the right side will say either "Up to date" or "Upgrade available".

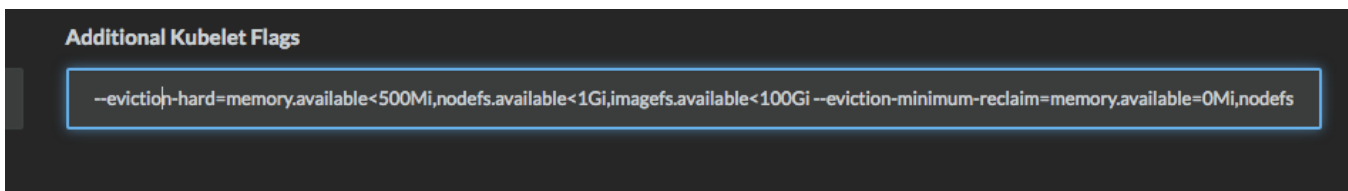


- Click the one beside the "+ kubernetes stack" and it will bring up a menu where you can select the environment template you want to tweak



- Find the field for Additional Kubelet Flags and input the desired amount of resources you want to reserve. I used the example from the K8s docs but you should adjust to what makes sense for your environment.

```
--eviction-hard=memory.available<500Mi,nodefs.available<1Gi,imagefs.available<5Gi
--eviction-minimum-reclaim=memory.available=0Mi,nodefs.available=500Mi,imagefs.available=2Gi
--system-reserved=memory=1.5Gi
```



- The stack will upgrade and it will basically add some startup parameters to the kubelet container and bounce them. Click on "Upgraded: Finish Upgrade" to complete the upgrade.

Stack: <span>⌵</span> kubernetos		Add Service <span>⌵</span>		Upgraded: Finish Upgrade		⌵ Upgraded	⌵
⚙ Active	addon-starter <span>⌵</span>	Image: rancher/k8sv1.10.5-rancher1-4		Service	1 Container	⌵	⌵
⚙ Active	controller-manager <span>⌵</span>	Image: rancher/k8sv1.10.5-rancher1-4		Service	1 Container	⌵	⌵
⚙ Started Once	etcd + 1 Sidekick <span>⌵</span>	Image: rancher/etcdv2.3.7-17		Service	6 Containers	⌵	⌵
⚙ Active	kubecti-shell <span>⌵</span>	Image: rancher/kubectidv0.8.7		Service	1 Container	⌵	⌵
⚙ Active	kubectid <span>⌵</span>	Image: rancher/kubectidv0.8.7		Service	1 Container	⌵	⌵
⌵ Upgraded	kubelet <span>⌵</span>	Image: rancher/k8sv1.10.5-rancher1-4		Service	12 Containers	⌵	⌵
⚙ Active	kubernetes + 1 Sidekick <span>⌵</span>	Image: rancher/k8sv1.10.5-rancher1-4		Service	2 Containers	⌵	⌵
⚙ Active	proxy <span>⌵</span>	Image: rancher/k8sv1.10.5-rancher1-4		Service	6 Containers	⌵	⌵
⚙ Active	rancher-ingress-controller <span>⌵</span>	Image: rancher/rb-service-rancherv0.9.4		Service	1 Container	⌵	⌵
⚙ Active	rancher-kubernetes-agent <span>⌵</span>	Image: rancher/kubernetes-agentv0.6.9		Service	1 Container	⌵	⌵
⚙ Active	rancher-kubernetes-auth <span>⌵</span>	Image: rancher/kubernetes-authv0.0.8		Service	1 Container	⌵	⌵
⚙ Active	scheduler <span>⌵</span>	Image: rancher/k8sv1.10.5-rancher1-4		Service	1 Container	⌵	⌵

- Validate the kubelet flags are now passed in as arguments to the container:

Service: ⌵ kubelet

In kubernetos

Type: Service

Scale: Global

Image: rancher/k8s:v1.10.5-rancher1-4

Entrypoint: None

Command: kubelet --kubeconfig=/etc/kubernetes/ssl/kubeconfig --register-node=true --cloud-provider=rancher --allow-privileged=true --healthz-bind-address=0.0.0.0 --cluster-dns=10.43.0.10 --fail-swap-on=false --cluster-domain=cluster.local --network-plugin=cni --cni-conf-dir=/etc/cni/managed.d --anonymous-auth=false --volume-plugin-dir=/var/lib/kubelet/volumeplugins --client-ca-file=/etc/kubernetes/ssl/ca.pem --pod-infra-container-image=gcr.io/google\_containers/pause-amd64:3.0 --tls-cipher-suites=TLS\_ECDHE\_ECDSA\_WITH\_AES\_128\_GCM\_SHA256,TLS\_ECDHE\_ECDSA\_WITH\_AES\_256\_GCM\_SHA384 --eviction-hard=memory.available<500Mi,nodefs.available<1Gi,imagefs.available<10Gi --eviction-minimum-reclaim=memory.available=0Mi,nodefs.available=500Mi,imagefs.available=0Mi --system-reserved=memory=1.5Gi

Ports Containers Labels Links Log

State ⌵	Name ⌵	IP Address ⌵	Host ⌵
○ Running	kubernetes-kubelet-10	None	k8s-1
○ Running	kubernetes-kubelet-5	None	k8s-3
○ Running	kubernetes-kubelet-6	None	k8s-4
○ Running	kubernetes-kubelet-7	None	k8s-5
○ Running	kubernetes-kubelet-8	None	k8s-10
○ Running	kubernetes-kubelet-9	None	k8s-2

- That's it!

Hopefully this was helpful and will save people time and effort reviving dead K8s worker nodes due to resource starvation from containers that do not have any resource limits in place.