# Live Migration of Production Workloads from Apache Mesos PaaS to Kubernetes

*Maria Camacho & Gufran Lutful, Nokia*

# Who we are

Maria Camacho          Gufran Lutful

*"A picture is worth a thousand words"*

**Nokia** has a comprehensive portfolio of network equipment, software, services and licensing opportunities across the globe for communications service providers.

With its commitment to innovation, driven by the award-winning Nokia Bell Labs, Nokia is a leader in the development and deployment of 5G networks.
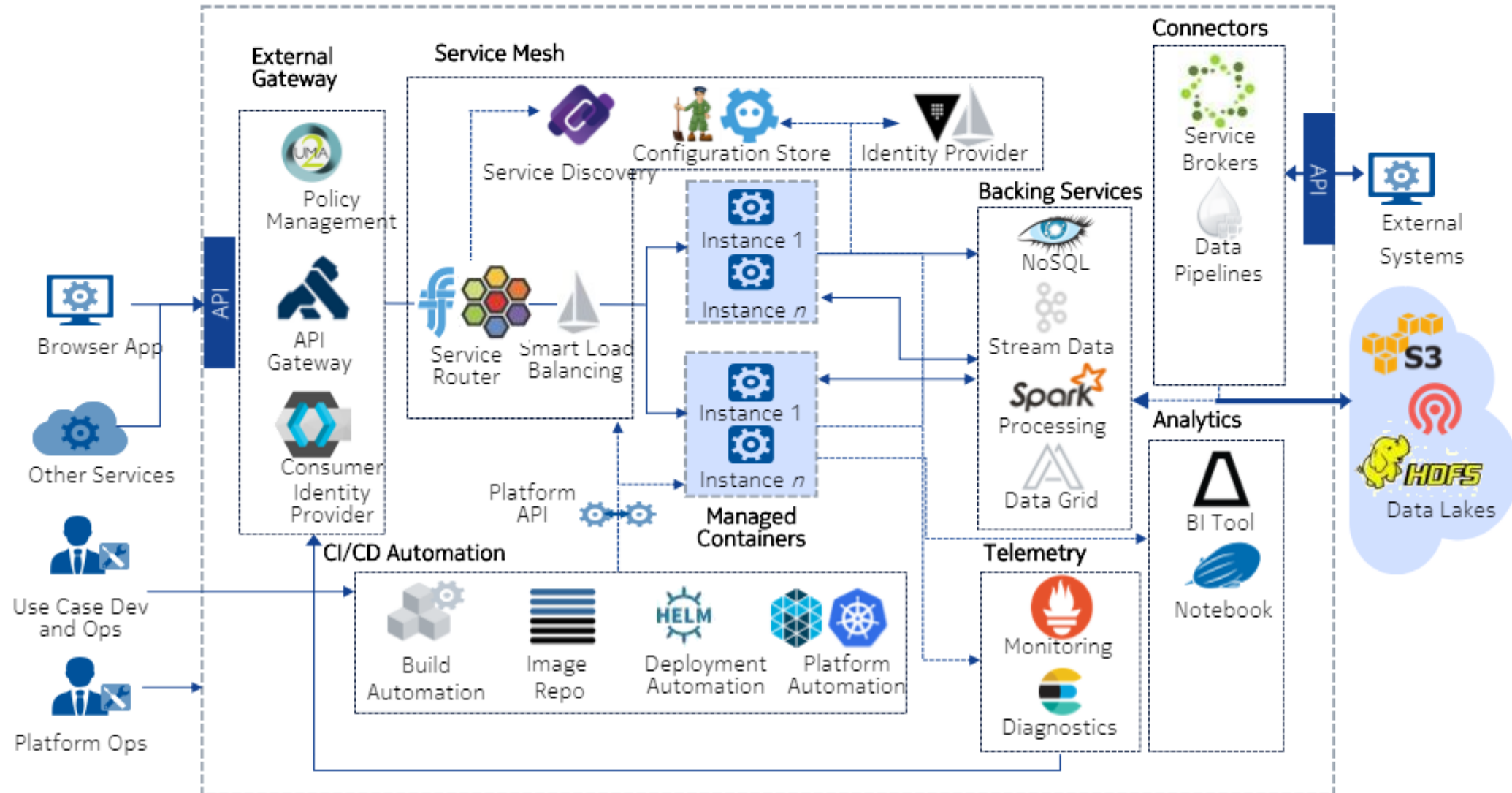
Nokia is still connecting people ;)

NOKIA

©priestmangoode

# About the project

# Adoption of Kubernetes
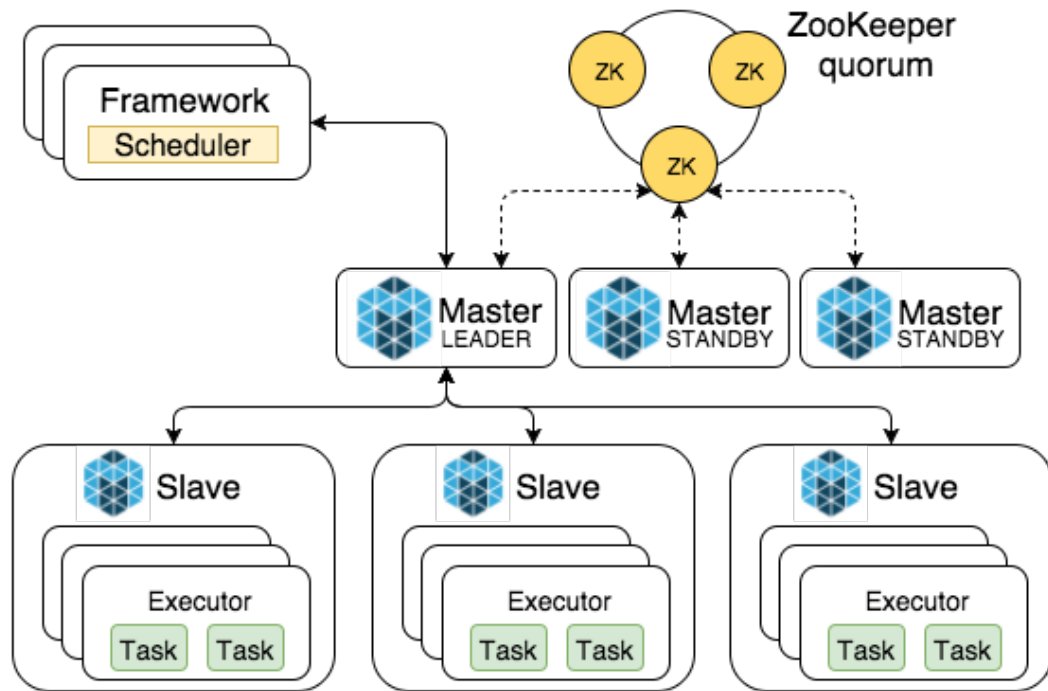
By 2018...

Kubernetes has become:

- ✓ Industry leading container orchestrator
- ✓ One of the top projects on GitHub: in a top position in stars, and No. 1 in terms of activity
- ✓ The centre of a growing community
- ✓ Quickly reaching production-level maturity

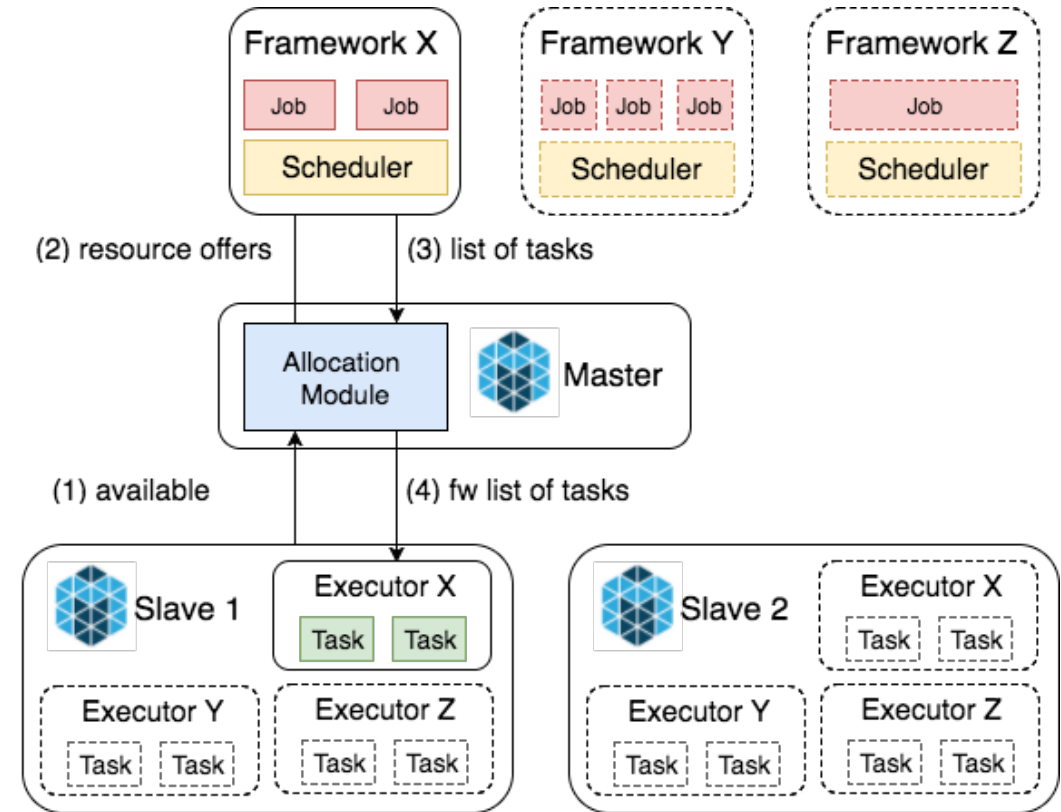But there were limitations too:

- ✓ Hard to run big data workloads with Apache Spark
- ✓ Not possible to seamlessly manage LPVs
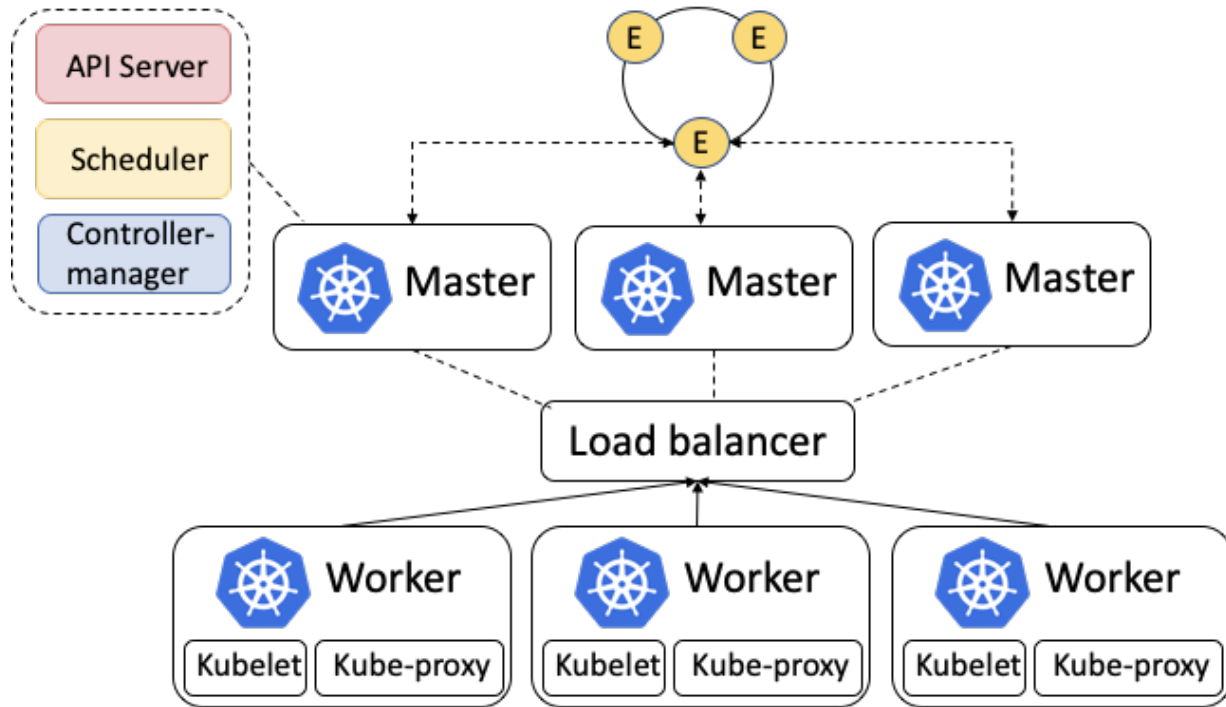
NOKIA

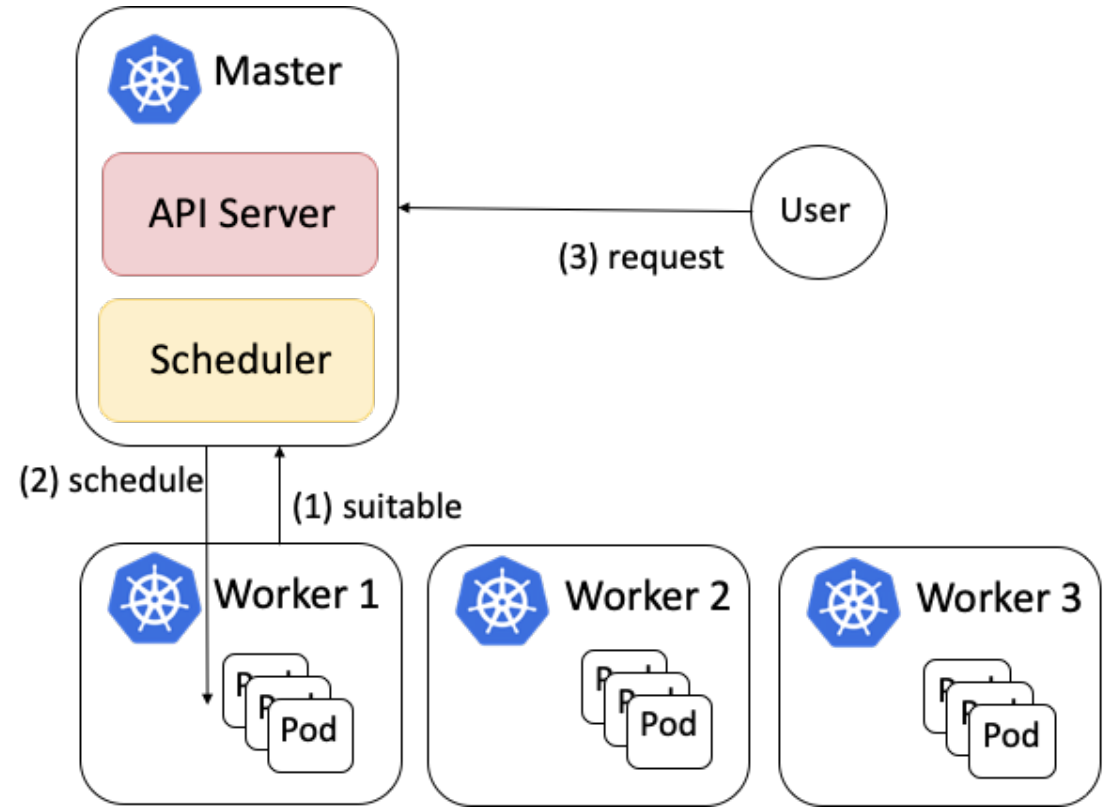# Mesos Overview

## Architecture
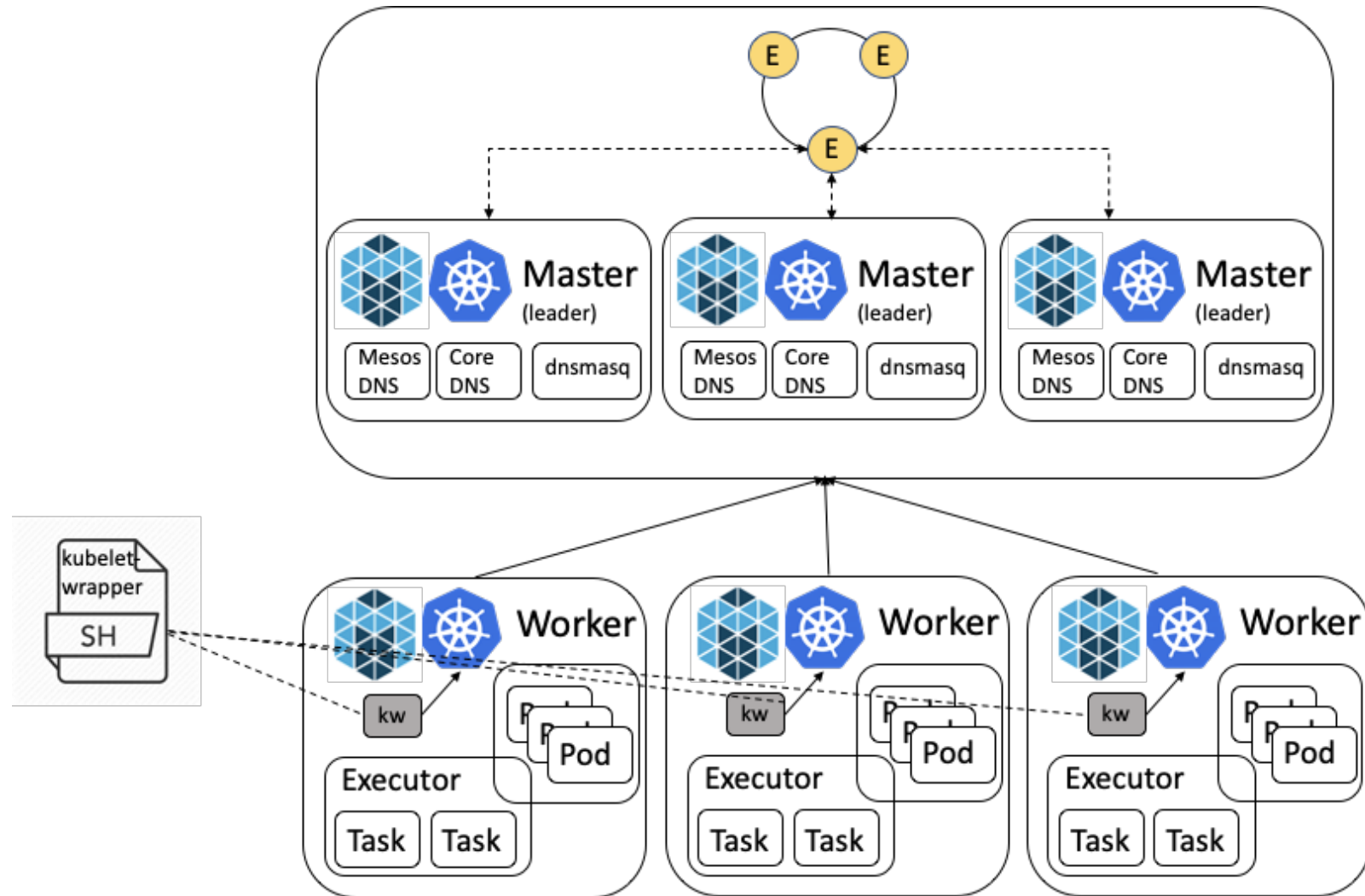


## Scheduling

# Kubernetes Overview

Architecture

Scheduling

# Mesos & Kubernetes Together

# Kubelet Wrapper in Marathon

# Kubelet Wrapper in VM-Image

| | |
|---|---|
| 📁 base | Remove Flocker from VM |
| 📁 cilium | fix policy filter trigger |
| 📁 dnsmasq | ESPOOBL-6048: Refactor LVM partitioning and fix resolvc... |
| 📁 docker | ESPOOBL-6657: Fix docker socket race issue |
| 📁 etcd3 | Update Etcd from 3.2.17 to 3.3.11 |
| 📁 flannel | ESPOOBL-5736: Install Cilium to VM images |
| 📁 health-checks | fix |
| 📁 kubernetes | ESPOOBL-8165, hotfix for removing kube-resource-alloca... |
| 📁 load-balancer | switch gitlabe1 to e2 for ava-core deps |
| 📁 marathon | Increase Marathon memory 1G -> 2G due to prod feedba... |
| 📁 mesos | ESPOOBL-7483: Add rootflags to enable quotas |
| 📁 nexus-preload/tasks | Use Artifactory proxy for docker images |
| 📁 node-config | ESPOOBL-8165, hotfix for enabling kube-resource-allocat... |
| 📁 openproxy | Introduce version 2.5.8 |
| 📁 rexray | add Rexray README doc |
| 📁 tests | ESPOOBL-8165, K8s resource enhancements. |
| 📁 zookeeper | ESPOOBL-7423: Increase the maximum limit of concurren... |

📄 roles/kubernetes/tasks/main.yml

```
24        dest: /etc/kubernetes/
25
26    - name: Copy kubelet-resource-wrapper
27      template:
28        src: kubelet-resource-wrapper
29        dest: /usr/local/bin/
30        mode: 0755
```

📄 roles/node-config/files/etc/init-k8s

```
142
143      # Join cluster
144      /usr/local/bin/kubelet-resource-wrapper
145
146      # Enable kubelet resource allocator service to start at boot
```

NOKIA

# Launching Kubelet Wrapper in Marathon



```
kubernetes-services.yml 741 Bytes          Edit   Web IDE   Replace   Delete

1   marathon:
2   - data:
3       id: "/k8s/services-resources"
4       instances: "{{ k8s_services_instances | default(1) | to_json_number }}"
5       cpus: "{{ k8s_services_cpu | default(2) | to_json_number }}"
6       mem: "{{ k8s_services_mem | default(4196) | to_json_number }}"
7       cmd: "echo cpu $MARATHON_APP_RESOURCE_CPUS, memory $MARATHON_APP_RESOURCE_MEM > /etc/kube-resources/services; /usr/local/bin/kubelet-resource-wrapper wait"
8       constraints: [["hostname", "UNIQUE"]]
9       healthChecks:
10      - gracePeriodSeconds: 60
11        intervalSeconds: 30
12        timeoutSeconds: 5
13        maxConsecutiveFailures: 0
14        path: "/healthz"
15        protocol: "MESOS_HTTP"
16        port: 10248
17      upgradeStrategy:
18        minimumHealthCapacity: 0
19        maximumOverCapacity: 0
```

```
74
75    .PHONY: metadata-apps
76    metadata-apps:
77            @scripts/apps-cli app install $(APP_INSTALL_PATTERN)
78
79    .PHONY: system-apps
80    system-apps:
```

📄 **app.yml** 738 Bytes 📋

```
1    name: monitoring/prometheus
2    version: 1.0.0
3    api_version: v1
4
5    description: "Systems monitoring and alerting toolkit"
6    helm:
7      app_name: prometheus
8      app_namespace: ...
9    app_type: helm
10
11    resources:
```

📄 **app.yml** 258 Bytes 📋

```
1
2    name: workspaces/couchdb
3    version: 2.3.0-1.1.0
4    api_version: v1
5
6    description: "CouchDB"
7
8    resources:
9      docker_images:
10       workspaces_couchdb:
11         repo: "https://gitlabe1.ext.net.nokia.com/
```

📄 src/ansible/apps-**metadata**_v1.yml

```
5    roles:
6
7    - role: metadata-deploy/read-app-config
8      run_once: true
9      tags: config
10
11   - role: metadata-deploy/export-app-resources
12     run_once: true
13     delegate_to: localhost
14     tags: config, export-resources
15
16   - role: metadata-deploy/export-network-policies
17     run_once: true
18     delegate_to: localhost
19     tags: config, export-resources
20
21   - role: metadata-deploy/helm-deploy
22     run_once: true
23     delegate_to: localhost
```

📄 src/ansible/apps-**metadata**_v1.yml

```
25       when: k8s_enabled | default(false)
26
27   - role: metadata-deploy/apply-network-policies
28     run_once: true
29     delegate_to: localhost
```

📄 src/ansible/apps-**metadata**_v1.yml

```
48   roles:
49
50   - role: metadata-deploy/deploy
51     tags: deploy
```

# Lessons learnt

**The strategy:**

- ✓ Spike to study the possible options of migration
- ✓ Follow KISS principle
- ✓ Less is more
- ✓ Favour a release-driven migration
- ✓ Have proper documentation/guidelines for dev teams
- ✓ Have a rollback strategy

**The implementation:**

- ✓ Metadata driven deployment
- ✓ K8s and Mesos can share same host resources
- ✓ Dimension your cluster properly, including system resources
- ✓ Use dedicated CIDRs for each orchestrator
- ✓ Kubelet can be run with no resources. Required for pod eviction

**The benefits:**

- ✓ Seamless sharing of resources between orchestrators
- ✓ Hosting selected workloads on each orchestrator
- ✓ Managing network traffic between orchestrators
- ✓ Internal DNS sharing
- ✓ Independent block storage management
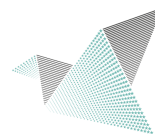
...and much more

NOKIA

# Bonus info

Some tools we used and loved: