# Scaling Up Logging and Metrics

Ricardo Lourenço < rlourenc@redhat.com >
Elvir Kuric < ekuric@redhat.com >

# Agenda

- Performance and Scale
- OCP Logging stack
  - Overview
  - Test harness
  - Large clusters
  - Scaling Up
  - 
- OCP Metrics stack
  - Overview

- Tools used
- Sample results

**OPEN**SHIFT
by Red Hat

# Performance and Scale

# We test individual components for limits. Then we scale them up.

We focus on OpenShift's control plane.

Then we target any underlying systems.

Components: ElasticSearch, Kibana, Cassandra, FluentD, Heapster, Java, Ruby...

# Logging stack overview

Users need a way to collect statistics about nodes in a cluster and send those logs to an external database.

- Elasticsearch is a distributed system by itself.
    - Collects logs from node services and all containers in the cluster.
    - Should be deployed with redundancy and persistent storage for scale and high availability.
    - Labels (region=infra) and Node Selectors.
- Fluentd
    - Gathers log entries from nodes, enriching them with metadata.
    - Feeds them into Elasticsearch.
    - Deployed using DaemonSets, Labels (logging-infra-fluentd).
- Kibana
    - presents a web UI for browsing and visualizing logs in Elasticsearch.

# MaxRequestsInFlight

```
servingInfo:
  bindAddress: 0.0.0.0:8443
  bindNetwork: tcp4
  certFile: master.server.crt
  clientCA: ca.crt
  keyFile: master.server.key
  maxRequestsInFlight: 500
  requestTimeoutSeconds: 3600
```

Partial list of consumers:

- FluentD pods talk to the apiserver to get metadata for the pods such as its namespace, pod name, host.

- openshift-node talks with the apiserver for command and control.

- SDN

```
$ kube-apiserver --help | grep inflight

        --max-requests-inflight int       The maximum number of requests in flight at a given

time. When the server exceeds this, it rejects requests. Zero for no limit. (default 400)
```

OPENSHIFT
by Red Hat

# ReplicationController

*The ReplicationController makes sure that a pod or homogeneous set of pods are always up and available. If there are too many pods, it will kill some. If there are too few, the ReplicationController will start more.*

- The Replication Controller is limited by the scheduler and scheduling cycle.
- The DaemonSet binds directly to the node, and bypasses the scheduler entirely.

Burst replicas on both daemonsets and rc is set to 500.

This essentially creates a step f(n) from 0-500 ~O(1) time.

```
apiVersion: v1
kind: ReplicationController
metadata:
  name: nginx
spec:
  replicas: 3
  selector:
    app: nginx
  template:
    metadata:
      name: nginx
      labels:
        app: nginx
    spec:
      containers:
      - name: nginx
        image: nginx
        ports:
        - containerPort:80
```

# DaemonSet

*A way to run a daemon on every node, or on a certain set of nodes in a cluster.*
*Use a a DS instead of a ReplicationController for pods that provide a machine-level function, such as machine monitoring or machine logging.*

Use cases:
- building a sharded datastore.
- running a logger on every node.
- cluster-level applications with strong node ties, and Openshift / Kubernetes node services.

The DaemonSet is used to run Fluentd on every node and send the data to a service like ElasticSearch for analysis.

pkg/controller/replication/replication_controller.go

```
// performance requirements for kubernetes 1.0.
BurstReplicas = 500

// The number of times we retry updating a replication
controller's status.
```

```
// It resumes normal action after observing the watch events for
them.
```

[pkg/controller/daemon/daemoncontroller.go](pkg/controller/daemon/daemoncontroller.go)

```go
// performance requirements for kubernetes 1.0.
BurstReplicas = 500

// If sending a status update to API server fails, we
retry a finite number of times.


// It resumes normal action after observing the watch
events for them.
```

# Limits

# ES

## Queue capacity

At 10250 pods (50 logger pods @ ~205 OCP nodes) Logging a 256 Byte string per second we've hit the ES bulk queue capacity.

# Fluentd

## Buffer chunk limit

The memory buffer plugin provides a fast buffer implementation. It uses memory to store buffer chunks. When Fluentd is shut down, buffered logs that can't be written quickly are deleted.

# Logging throughput test

- Features:
  - Full logging project deployment.
  - Automated testing with pbench integration.
    - Loading up a cluster with logging traffic at different rates.
    - Collect pbench data while containers are logging.
  - Other logging drivers are also supported through the -d parameter.
    - Json-file, fluentd, journald

# Test description

**Phase 1 - Before the loggers and pbench start**

Delete all indexed data so that data from previous runs doesn't add up.
   $ oc exec $es_pod -- curl -XDELETE "$es_url/*"

Gets ES node stats through the API
   $ oc exec $es_pod -- curl $es_url/_stats?pretty

Captures disk usage under /var/lib/origin/openshift.local.volumes/pods/

**Phase 2 - Generate load**

Start X logger pods across N OCP cluster nodes, logging pre-defined random string at a chosen logging rate.
Registers pbench tools across chosen cluster nodes: iostat mpstat pidstat proc-vmstat sar turbostat

**Phase 3 - After the run finishes.**

Optimize and flush indexes.
   $ oc exec $es_pod -- curl -XPOST "$es_url/_refresh"
   $ oc exec $es_pod -- curl -XPOST "$es_url/_optimize?only_expunge_deletes=true&flush=true"

Like in phase 1, gets ES, OC stats, logs and disk usage for your pbench --remote nodes.

# Sample output

# 10n-no_limits-2048ll-1000wait-300s

# 30wppn - 205nodes - 1hour test.
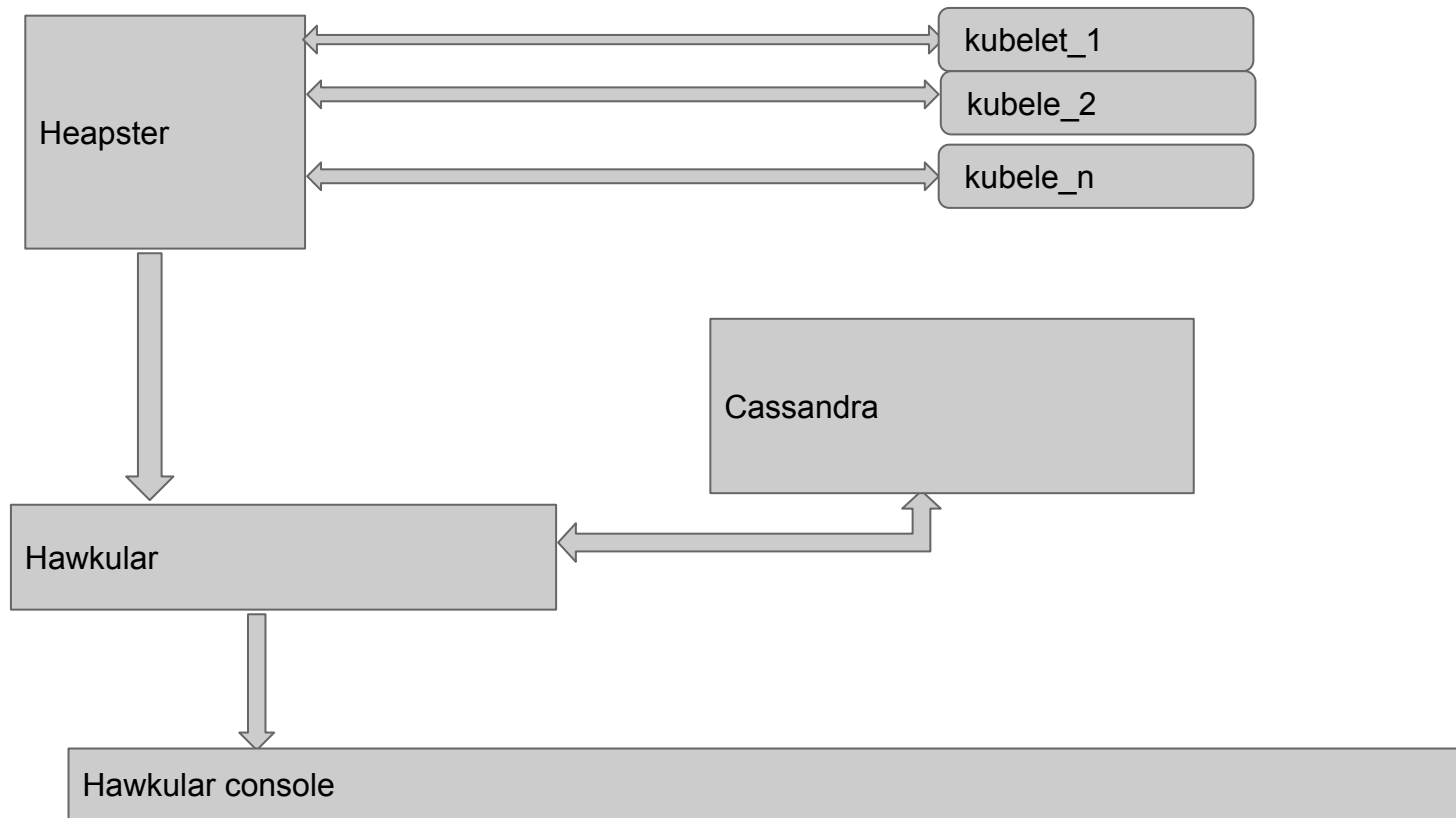
# Recommendations

- It is highly recommended to pre-pull the fluentd images (> 200 nodes).
  - ImagePullPolicy: IfNotPresent should be set in the deployer.yaml file.
- After deploying the region=infra pods, labeling your nodes should be done in steps/batches of 20 nodes at a time while keeping an eye open for any "MatchNodeSelector|CrashLoopBackOff|ImagePullFailed|Error" status.
- Use labels and selectors to differentiate infrastructure and application nodes
- Use large persistent volumes (PVs) for infrastructure nodes running logging (Elasticsearch) and metrics (Cassandra)

# OpenShift Metrics

# OCP Metrics stack - Overview

- Cassandra
  - Serves as datastore  - all metrics data are written to Cassandra db

- Hawkular metrics storage  engine

- Heapster
  - Collects metrics data from pods

Openshift metrics high level overview

# OpenShift Metrics console

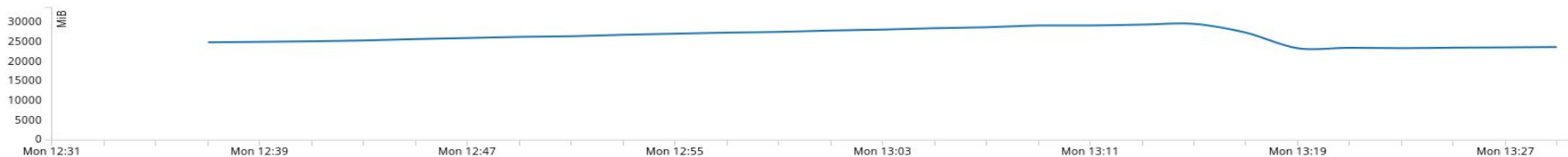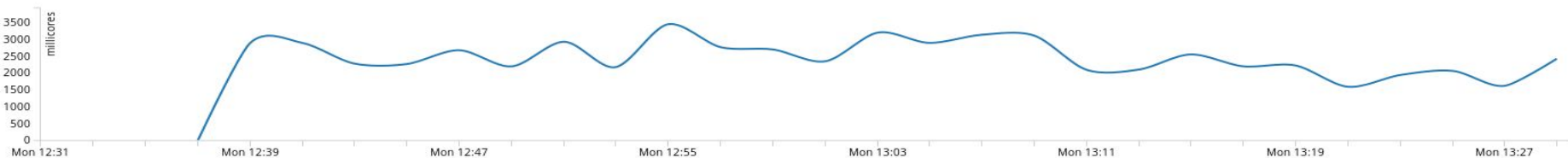# OCP Metrics configuration

- All metrics components are are delivered as images and starting deploying metrics is easy as starting metrics pods

- metrics.yaml  ( from upstream : https://github.com/openshift/origin-metrics )

OPENSHIFT
by Red Hat

# cassandra pod

- Pods ( = cassandra nodes ) form Cassandra cluster
- USE_PERSISTANT_STORAGE=True
    - PV will be mounted in cassandra pod /cassandra_data
- Persistent storage ensures data persistence
- Possible to scale out to more Cassandra pods (=nodes)

# Cassandra pod

- USE_PERSISTANT_STORAGE=false
- No data persistence - if pod dies data will be lost
- Watch /var/lib/origin/ if persistent storage not used - can fill it over time

# Heapster metrics component

- runs as Pod
- Heapster gathers metrics data from OCP cluster
- It gets metrics for every pod - across all namespaces
- Send these data to Hawkular metrics via API

# Hawkular

*"It provides means of creating, accessing and managing historically stored metrics via an easy to use json based REST interface."*

*Per :* [https://github.com/openshift/origin-metrics](https://github.com/openshift/origin-metrics)

OPENSHIFT
by Red Hat

# OCP metrics setup

- Manual setup
  - Enter commands on master
  - Edit /etc/master/master-config.yaml
  - Restart master
- Using advanced installer
  - Set up all in ansible playbook

https://docs.openshift.com/enterprise/3.1/install_config/cluster_metrics.html

OPENSHIFT
by Red Hat

# Scaling metrics pods

- We can add more
  - Cassandra pods
  - Hawkular pods
- Only one heapster pod per cluster

OPENSHIFT
by Red Hat

# OpenShift metrics

- Use persistent storage for data - always

    – OpenShift metrics supports dynamic storage provisioning

    – DYNAMICALLY_PROVISION_STORAGE=true

- Run metrics pods on 'infra' nodes

OPENSHIFT
by Red Hat

# OpenShift metrics configuration tips

- Watch system usage where metrics pods are running
- One set of metrics pods can handle ~10k pods - just guidelines

OPENSHIFT
by Red Hat

# OpenShift metrics configuration tips

- For monitoring more pods, necessary to scale out metrics pods

- Cassandra is not best candidate for network storage, eg. NFS

  - [http://docs.datastax.com/en/archived/cassandra/1.2/cassandra/architecture/architecturePlanningAntiPatterns_c.html](http://docs.datastax.com/en/archived/cassandra/1.2/cassandra/architecture/architecturePlanningAntiPatterns_c.html)

**OPENSHIFT**
by Red Hat

# OpenShift metrics scale numbers

- It was successfully ran on clusters up to 1000 nodes ( 210 and 981 nodes cluster )
- Cassandra storage requirements are
  - 2.5 GB / day  and monitored 1000 pods
  - 17.5 GB / day and 10k pods
- METRICS_DURATION=7(days)
- METRICS_RESOLUTION=15(s)

**OPEN**SHIFT
by Red Hat

# OpenShift Metrics scaling issues

- Heapster is not catching up with massive number of pods monitored > 12k pods  in  OCP cluster
- Up to 12k pods it works good
- BZ 1388815
- **Monitoring Microservices on OpenShift with HOSA**

Hawkular project :

OPENSHIFT
by Red Hat

# Thing to follow : HOSA

- Monitoring Microservices on OpenShift with HOSA
  - Hawkular OpenShift Agent (HOSA)

Blog post from hawkular team :

http://www.hawkular.org/blog/2017/01/17/obst-hosa.html?_lrsc=4c8f78df-8f03-4630-bd54-128a24998599&sc_cid=7016000000011zELAAY

# Thank you!
# Q & A ?

GitHub

- https://github.com/redhat-performance
- http://github.com/openshift/svt

Ricardo Lourenço rlourenc@redhat.com; Elvir Kuric ekuric@redhat.com

**OPEN**SHIFT
by Red Hat