



*Distributed data stores
on Kubernetes
* and other things*

Alena Hall

@lenadroid

Prerequisites

- *Kubernetes cluster*
- *Basic understanding of Cassandra*

Prerequisites

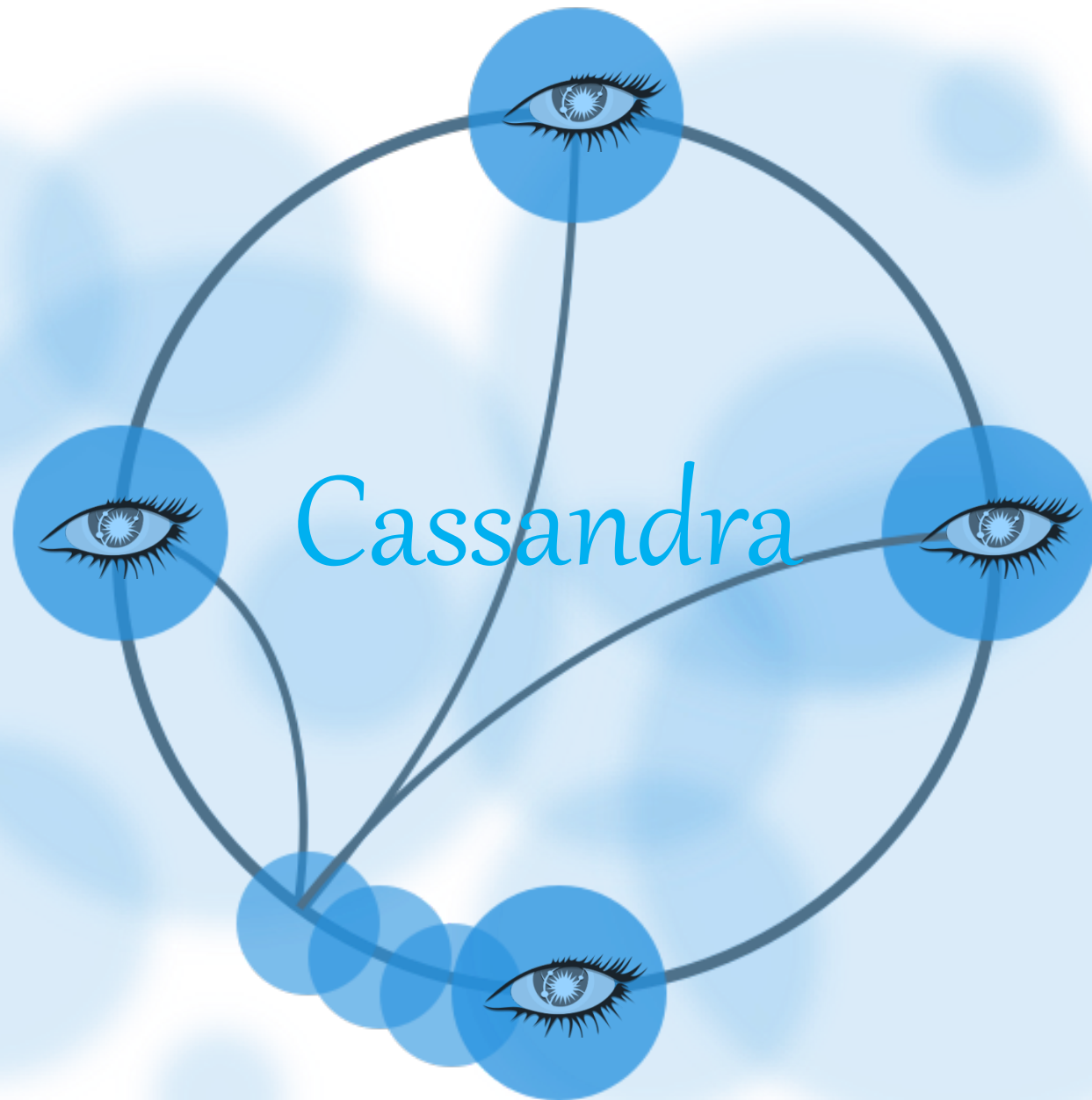
- Kubernetes cluster
- Basic understanding of Cassandra
- **Courage**

A horizontal, textured yellow brushstroke with irregular, feathered edges, serving as a background for the title text.

Motivation

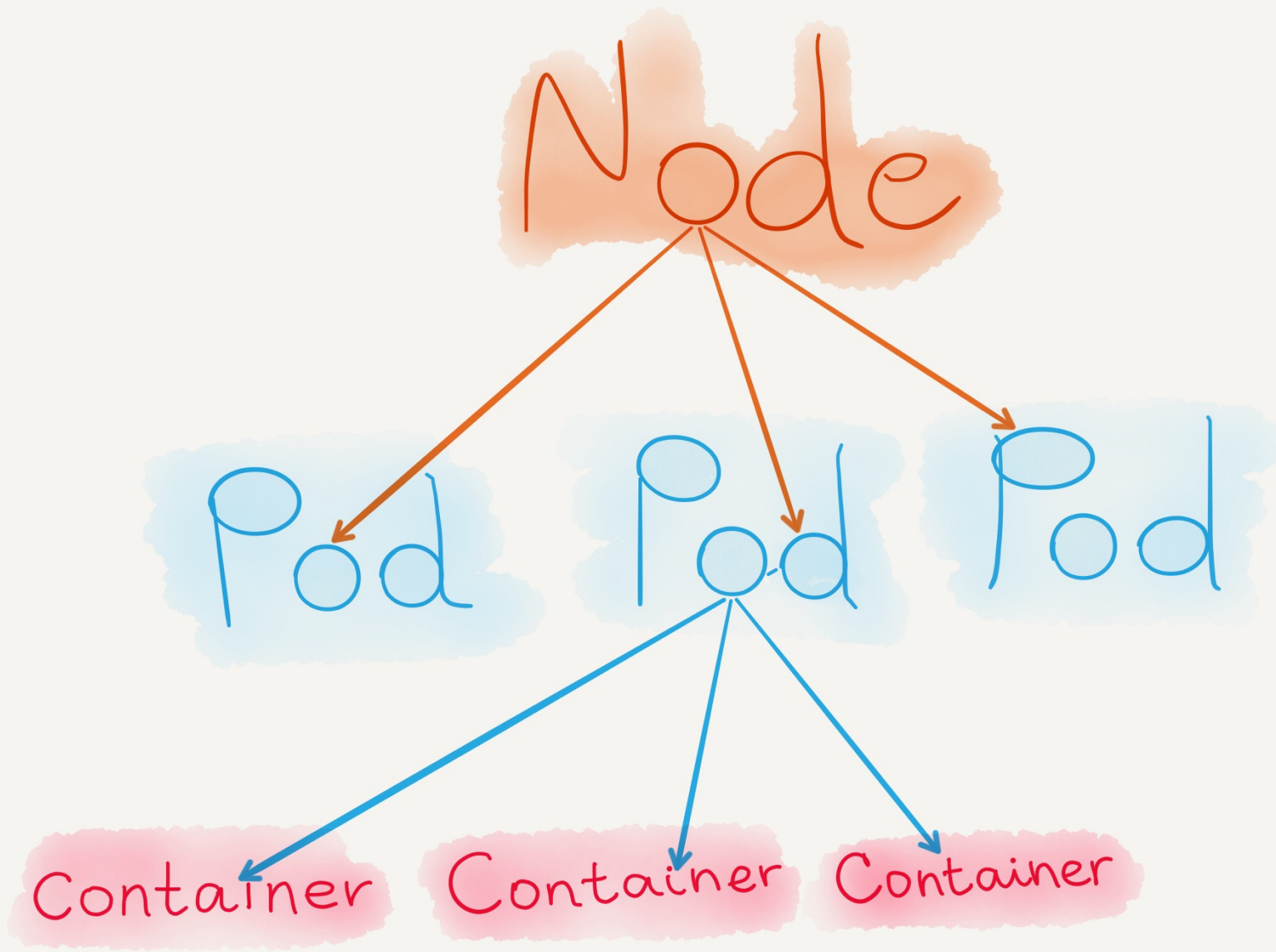
Talk plan

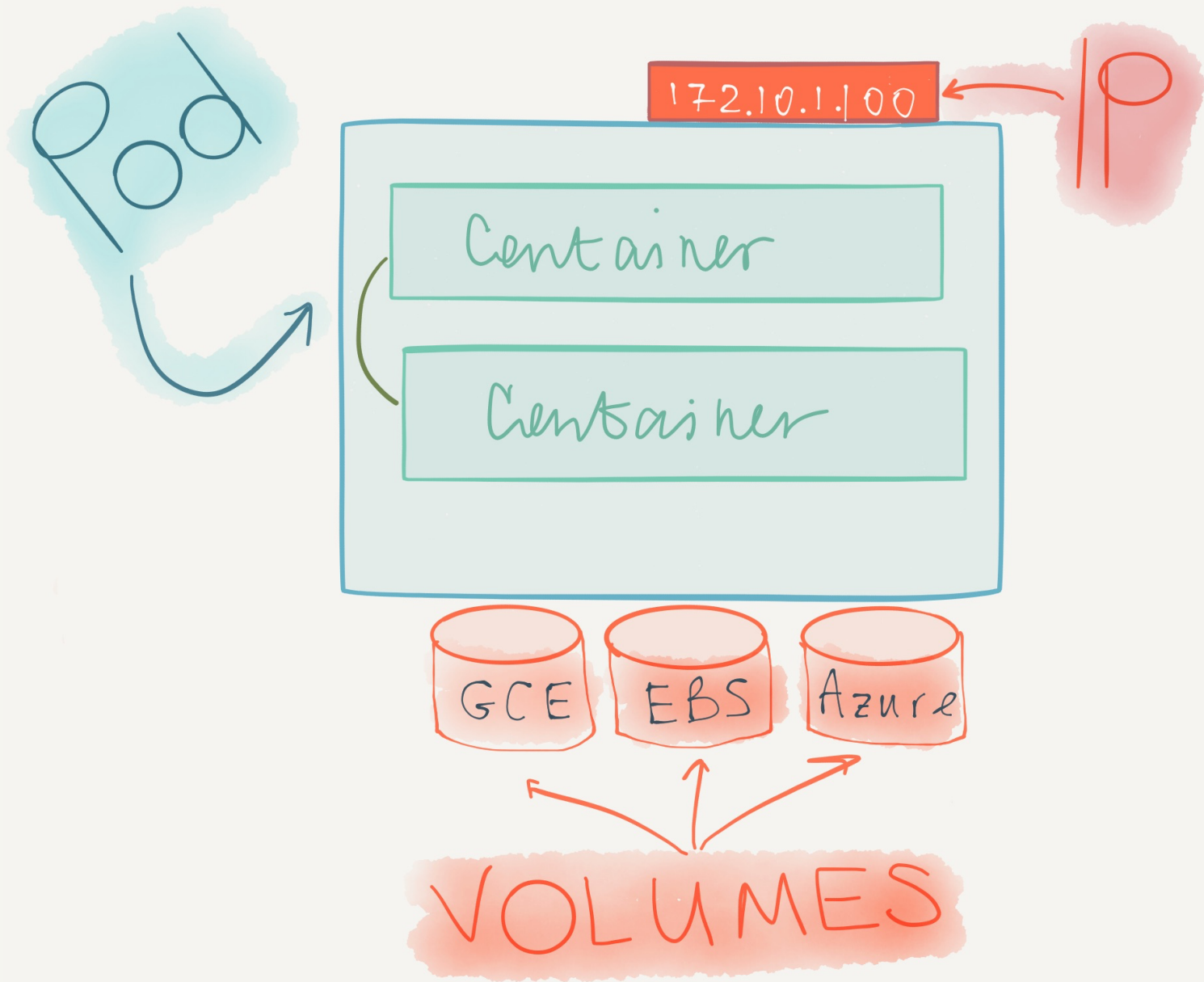
- ✓ Kubernetes and Cassandra refresher
- ✓ Stateful Sets, PVs, PVCs, Storage Classes, and more
- ✓ *Example*: Cassandra Stateful Set
- ✓ Kubernetes Jobs
- ✓ *Example*: Writing data into Cassandra using Jobs
- ✓ Spark jobs on Kubernetes
- ✓ *Example*: Spark and Cassandra on Kubernetes



- Distributed database, has a *ring topology*, uses *consistent hashing*
- Highly available system with *tunable consistency* and *hinted hand off*
- Data is grouped by the *partition key* and stored together physically
- Stores data in *SSTable* on disk and RAM and tracks operations in the *commit log*
- Data can be ordered by a *clustering key*
- Supports equality, inequality, aggregation queries and custom functions
- Hash ranges can be mapped to *Simple Nodes* or to *Vnodes*

Kubernetes Refresher







Art by @ashleymcnamara

@lenadroid

Kubernetes Deployments/Replica Sets

Keeps specified number of pod replicas running at any given time

Ephemeral containers



*Distributed databases - systems that require
stable persistence*

Kubernetes Volumes

What if we try...

Replicated containers with volumes?

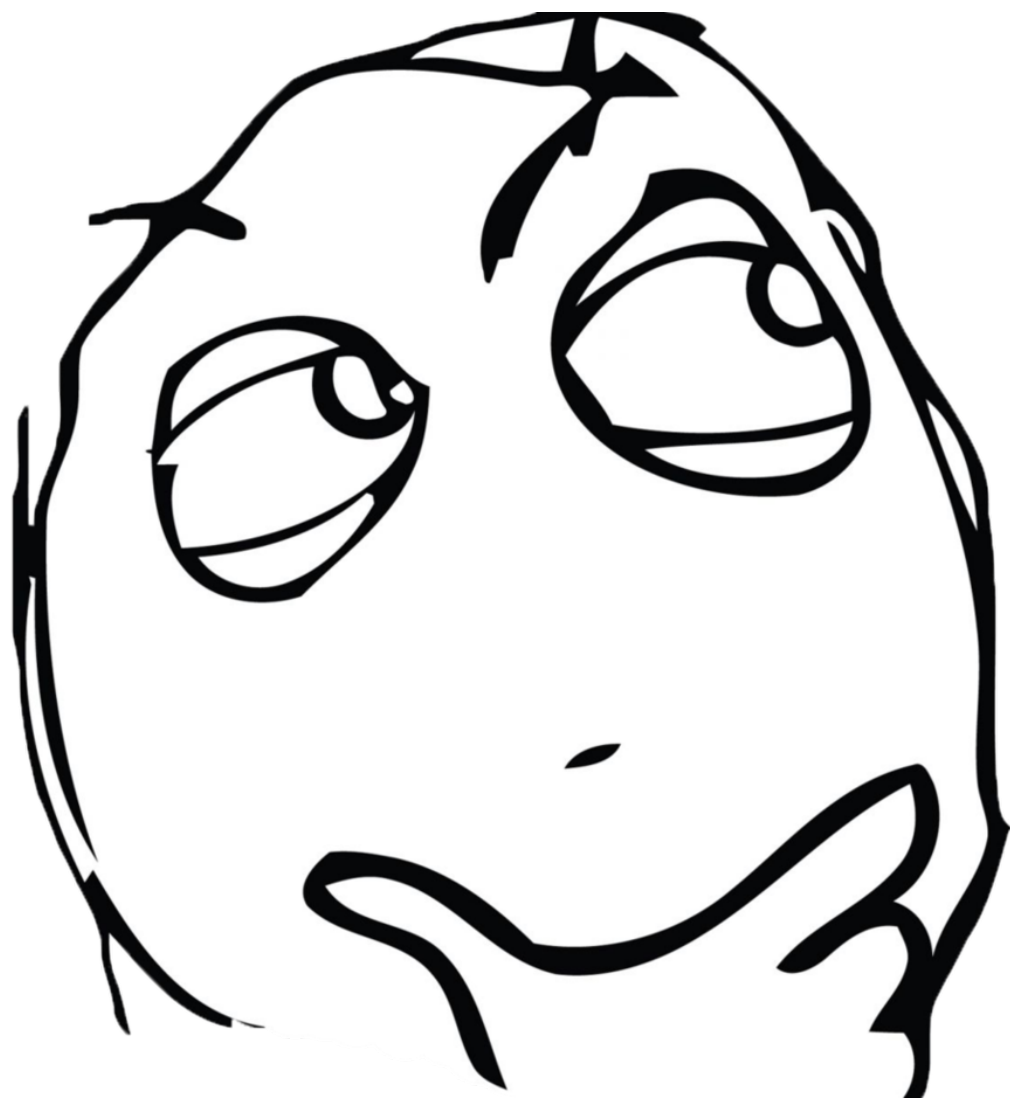
Standard containers are not persistent

Even combined with Kubernetes volumes, there are issues:

- Non-deterministic, **random** names
- Random order in which pods start, scale and terminate

So, replica sets with volumes won't work

@lenadroid



@lenadroid

Prior to creating the database cluster :

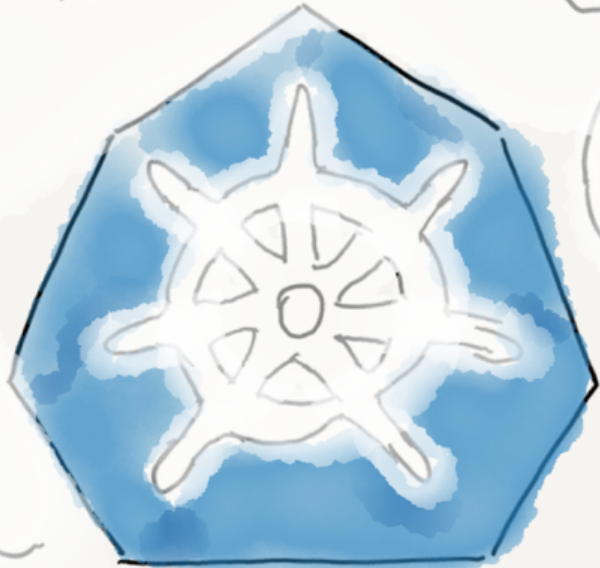
- We need to have certain **guarantees** about cluster nodes
- Database nodes should have **discoverable names**
- Nodes need to start in **predictable order**
- Stable persistent storage

A horizontal, textured orange brushstroke with irregular, feathered edges, serving as a background for the title text.

Stateful Sets

Stateful Set pods have identity ...

- Stable, **unique** network identifiers
- **Stable**, persistent storage, and link from pod to storage
- **Ordered** deployment, scaling, termination, rolling updates

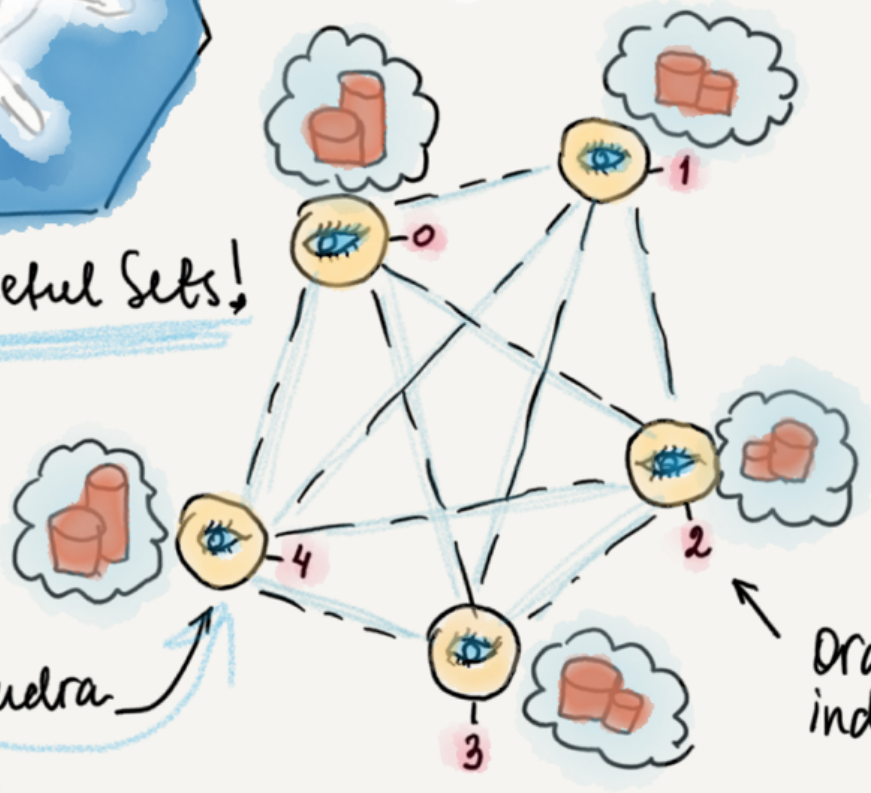


Kubernetes! Stateful Sets!

persistent volumes

Cassandra Pods

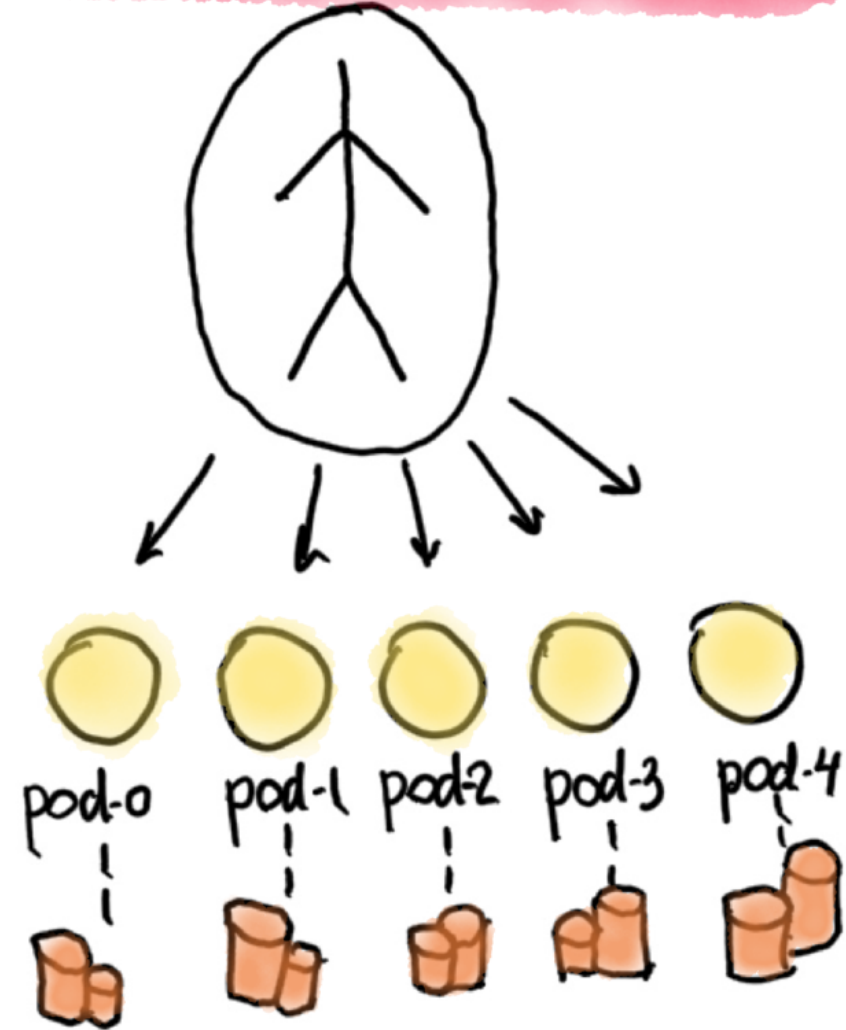
Ordinal index!



A horizontal, textured orange brushstroke with irregular, feathered edges, serving as a background for the title text.

Headless Service

Headless Service





! cassandra-service.yaml x



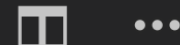
```
1  apiVersion: v1
2  kind: Service
3  metadata:
4    labels:
5      app: cassandra
6      name: cassandra
7  spec:
8    clusterIP: None
9    ports:
10   - port: 9042
11   selector:
12     app: cassandra
13
```





Storage Classes

! storage-class.yaml ×



```
1  {
2    "apiVersion": "storage.k8s.io/v1beta1",
3    "kind": "StorageClass",
4    "metadata": {
5      "annotations": {},
6      "labels": {
7        "kubernetes.io/cluster-service": "true"
8      },
9      "name": "managed-premium",
10     "namespace": ""
11   },
12   "parameters": {
13     "kind": "Managed",
14     "storageaccounttype": "Premium_LRS"
15   },
16   "provisioner": "kubernetes.io/azure-disk"
17 }
```






Persistent Volume Claims

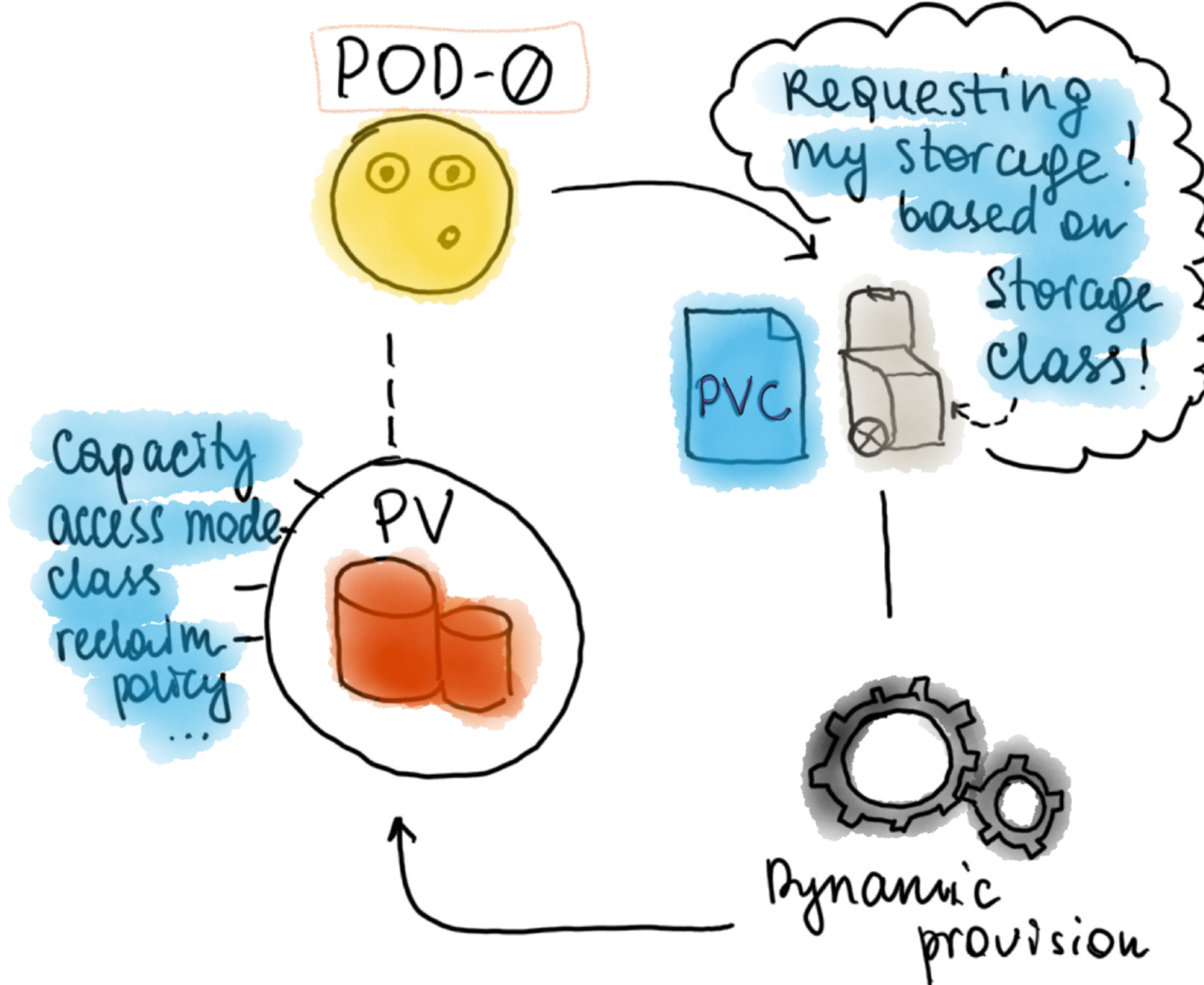


Persistent Volumes



*A story of persistent
volumes and persistent
volume claims*

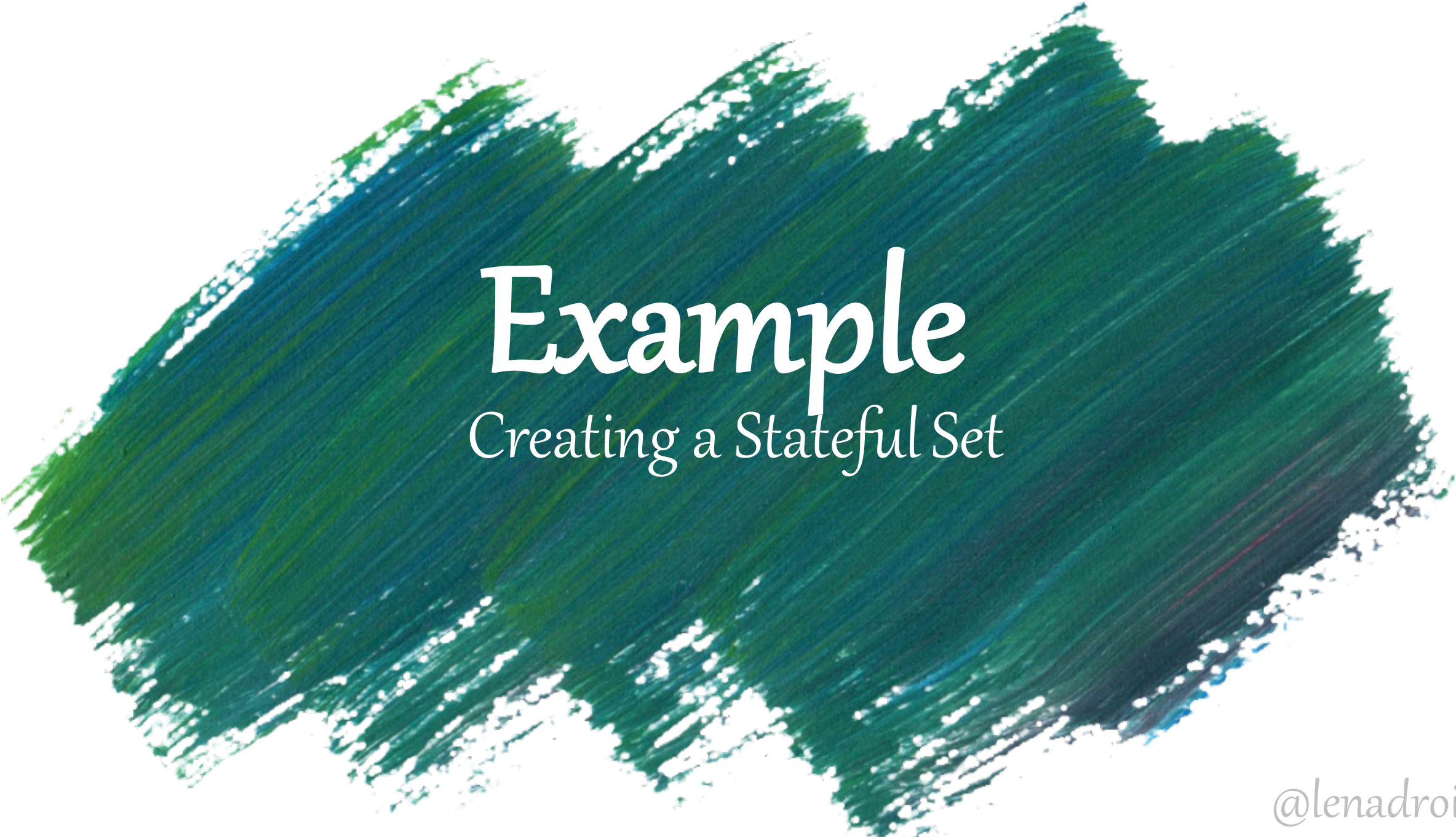
@lenadroid




```
! cassandra-statefulset.yaml x ...
1  apiVersion: apps/v1
2  kind: StatefulSet
3  metadata:
4    name: cassandra
5    labels:
6      app: cassandra
7  spec:
8    serviceName: cassandra
9    replicas: 5
10   selector:
11     matchLabels:
12       app: cassandra
13   template:
14     metadata:
15       labels:
16         app: cassandra
17     spec:
18       terminationGracePeriodSeconds:
19         1800
20       containers:
21         - name: cassandra
22           image:
23             gcr.io/google-samples/cassandra:v13
24           imagePullPolicy: Always
25           ports:
26             - containerPort: 7000
27               name: intra-node
28             - containerPort: 7001
29               name: tls-intra-node
30             - containerPort: 7199
31               name: jmx
32             - containerPort: 9042
33               name: cql
34           resources:
```

```
! cassandra-statefulset.yaml x ...
32   resources:
33     limits:
34       cpu: "2"
35       memory: 3.5Gi
36     requests:
37       cpu: "2"
38       memory: 3.5Gi
39   securityContext:
40     capabilities:
41       add:
42         - IPC_LOCK
43   lifecycle:
44     preStop:
45       exec:
46         command:
47           - /bin/sh
48           - -c
49           - nodetool drain
50   env:
51     - name: MAX_HEAP_SIZE
52       value: 512M
53     - name: HEAP_NEWSIZE
54       value: 100M
55     - name: CASSANDRA_SEEDS
56       value:
57         "cassandra-0.cassandra.
58         default.svc.cluster.
59         local"
60     - name:
61       CASSANDRA_CLUSTER_NAME
62       value: "chicago"
63     - name: CASSANDRA_DC
64       value: "DC1-chicago"
65     - name: CASSANDRA_RACK
66       value: "Rack1-chicago"
```

```
! cassandra-statefulset.yaml x ...
63     - name: POD_IP
64       valueFrom:
65         fieldRef:
66           fieldPath:
67             status.podIP
68   readinessProbe:
69     exec:
70       command:
71         - /bin/bash
72         - -c
73         - /ready-probe.sh
74     initialDelaySeconds: 15
75     timeoutSeconds: 5
76   volumeMounts:
77     - name: cassandra-data
78       mountPath: /cassandra_data
79   volumeClaimTemplates:
80     - metadata:
81         name: cassandra-data
82     spec:
83       accessModes: [
84         "ReadWriteOnce" ]
85       storageClassName:
86         managed-premium
87       resources:
88         requests:
89           storage: 1Gi
```



Example

Creating a Stateful Set



lena:gotochicago lenok\$ █

|

```
lena:gotochicago lenok$ kubectl get services
```

NAME	TYPE	CLUSTER-IP	EXTERNAL-IP	PORT(S)	AGE
cassandra	ClusterIP	None	<none>	9042/TCP	45m
cassandra-external	LoadBalancer	10.0.186.193	52.234.227.80	9042:30734/TCP	1d
kubernetes	ClusterIP	10.0.0.1	<none>	443/TCP	1d

```
lena:gotochicago lenok$ █
```



Example

Test for writing data

A horizontal, textured orange brushstroke with irregular, feathered edges, serving as a background for the title text.

Kubernetes Jobs

@lenadroid



! fsharp-job.yaml ×



```
1  apiVersion: batch/v1
2  kind: Job
3  metadata:
4    name: "process-item-$START"
5    labels:
6      jobgroup: fsharpjob
7  spec:
8    template:
9      metadata:
10     name: fsharpjob
11     labels:
12       jobgroup: fsharpjob
13     spec:
14       containers:
15         - name: fsharpjob
16           image: lenadroid/fsharp-job
17           command:
18             - fsharpi
19             - HousingData.fsx
20             - "$START"
21             - "$STEP"
22             - "$INCR"
23           resources:
24             requests:
25               memory: "1Gi"
26               cpu: "870m"
27             limits:
28               memory: "1Gi"
29               cpu: "870m"
30           restartPolicy: Never
31
```




prepare-jobs.sh ×



```
1  #!/bin/bash
2
3  jobCount=10
4  increment=1000
5  step=10000
6
7  jobDir=jobs
8
9  if [ -d "$jobDir" ]; then rm -Rf $jobDir; fi
10 mkdir $jobDir
11
12 for ((i=0; i <= $jobCount-1; i++))
13 do
14     startIndex=$((i * $step))
15     echo "Creating a job for rows starting from $startIndex"
16     cat fsharp-job.yaml | sed -e "s/\$START/$startIndex/" -e "s/\$STEP/$step/" -e "s/\$INCR/$increment/" >
17     ./$jobDir/job-$startIndex.yaml
18 done
```



1





Example

Creating Kubernetes Jobs

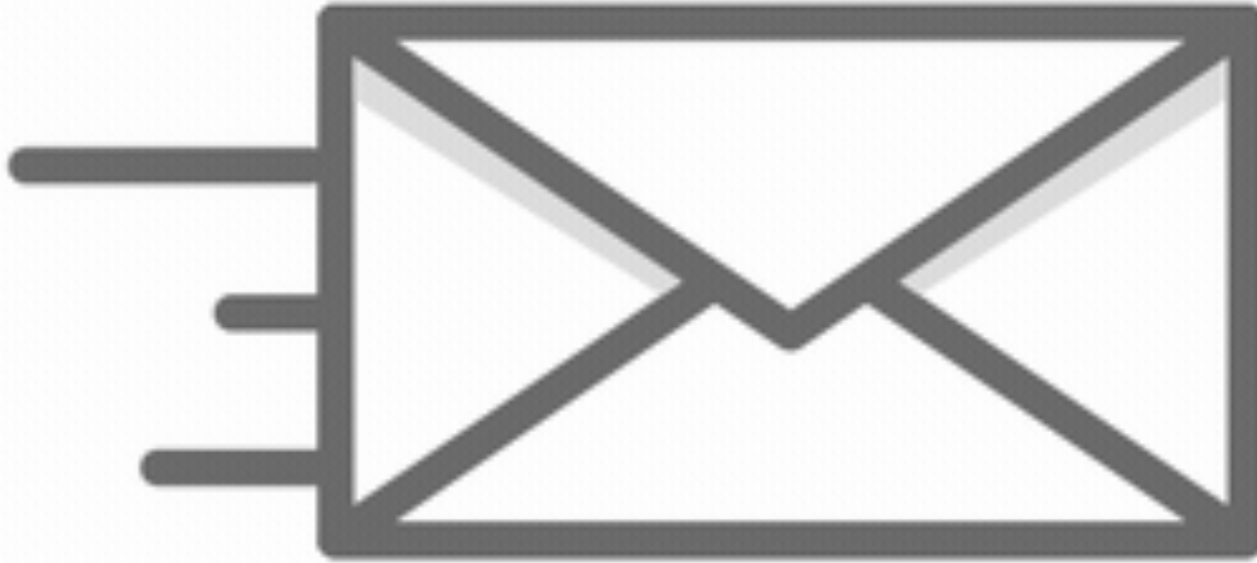
@lenadroid



1. bash

lena:gotochicago lenok\$ █

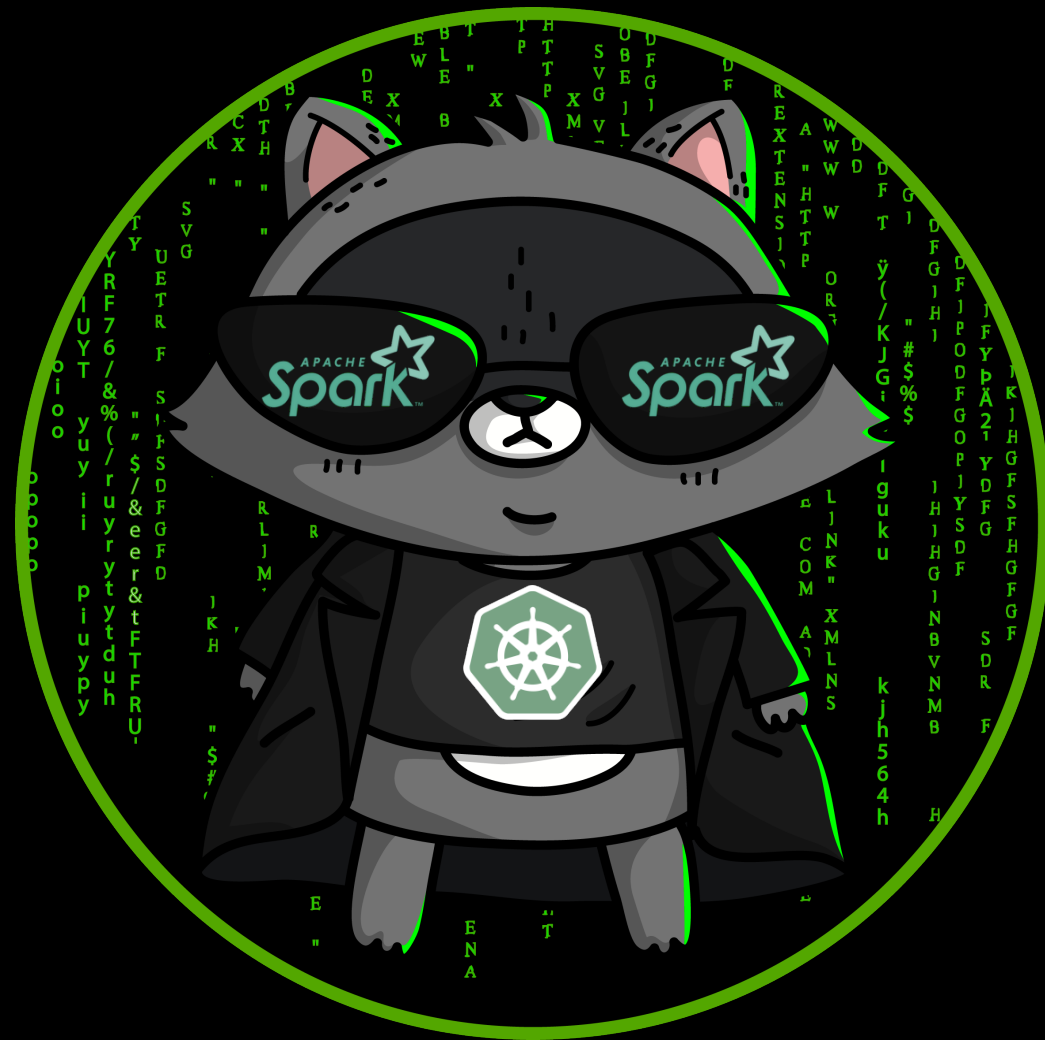
github.com/lenadroid/fsharp-job-kube



Things to note

- Persistent Volumes are not deleted automatically
- Stateful Sets were in beta before Kubernetes 1.9
- Node-affinity

Running Spark jobs with Kubernetes scheduler



Art by @ashleymcnamara

@lenadroid

Native Kubernetes support is very new for Spark!



Spark Driver



Spark Executors





Example

Running Spark jobs on Kubernetes
& with Cassandra

@lenadroid

lena:~ lenok\$

I

Is running Spark jobs on Kubernetes a good idea?

Maybe

- > For a test environment
- > To try out a custom version of Spark quickly
- > If you already have Kubernetes & want to use spare cycles
- > Don't want to have too many different orchestrators

Should I run it in production?

Not yet




BUT IT WILL GET BETTER

¹
YOU CAN HELP MAKE IT BETTER

Note:

*This is not the same as running standalone
Spark on Kubernetes*

These are -native spark jobs

github.com/lenadroid/goto-cassandra-spark



Resources

- *Cassandra blog post:*

<https://lenadroid.github.io/posts/cassandra-docker-fsharp.html>

- *Distributed systems algorithms talk:*

<https://www.safaribooksonline.com/library/view/oscon-2017-/9781491976227/video306675.html>

- *Stateful Sets article:*

<https://lenadroid.github.io/posts/stateful-sets-kubernetes-azure.html>

- *All demos on GitHub @ lenadroid*

- *Twitter @ lenadroid*



@lenadroid

Alena Hall



- Senior Cloud Developer Advocate, Engineer @ Azure
- Member of F# Software Foundation Board of Trustees
- Lives in beautiful Seattle, WA
- Loves doing data science, machine learning, functional programming and distributed systems at scale

  @lenadroid

Thank you!

Alena Hall   *@lenadroid*