

So You Want To Run a Data-Intensive System On Kubernetes

@lenadroid

#SeattleScalability

Agenda

Distributed systems for big data

Running them on Kubernetes?

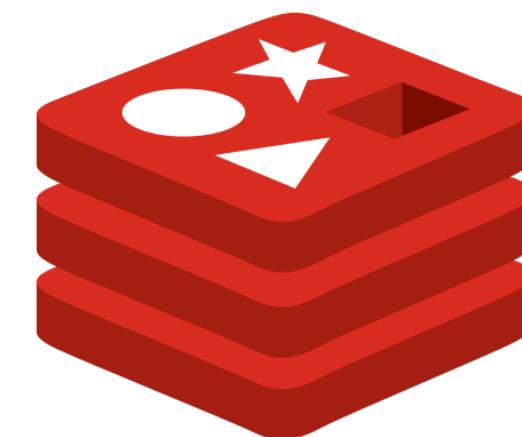
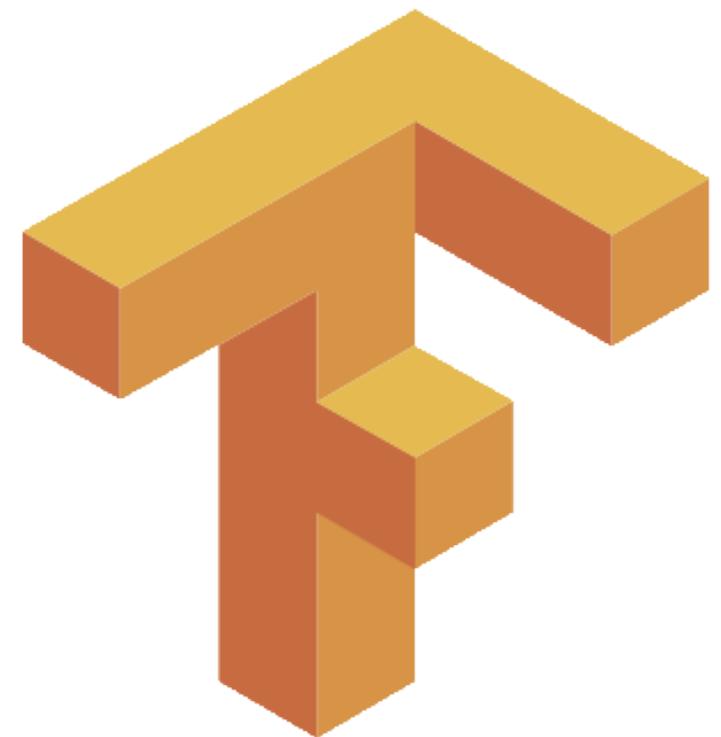
Getting them to “work”

Challenges

Work in progress and solutions

Future

Summary



redis

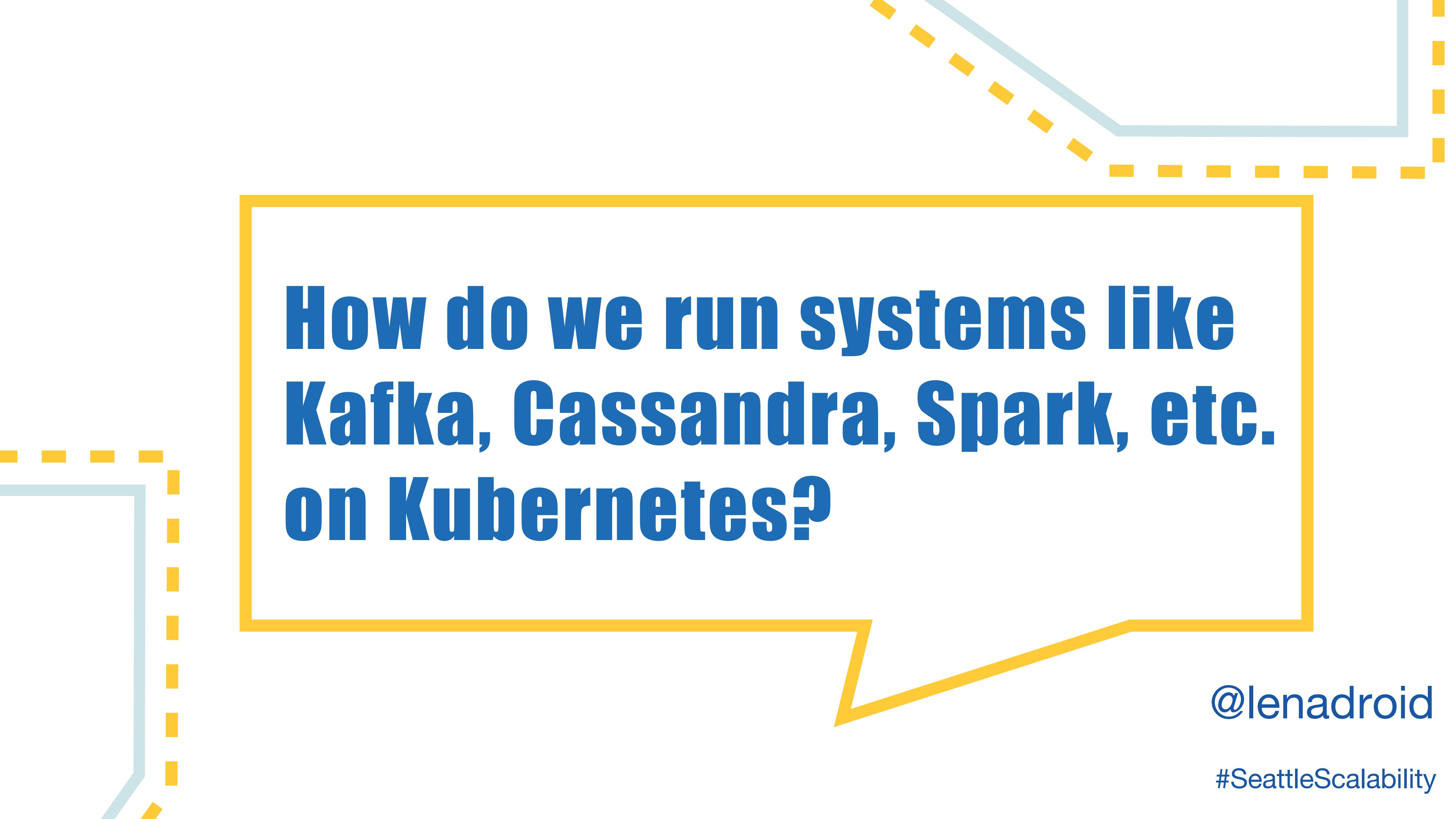
@lenadroid
#SeattleScalability

Data systems on Kubernetes



@lenadroid

#SeattleScalability



How do we run systems like Kafka, Cassandra, Spark, etc. on Kubernetes?

@lenadroid

#SeattleScalability

Interpretation #1

**What Kubernetes abstractions
can we use to run these
distributed systems?**

@lenadroid

#SeattleScalability

Interpretation #2

**What are the steps to simply
get those distributed systems
up and running based on these
Kubernetes abstractions?**

@lenadroid

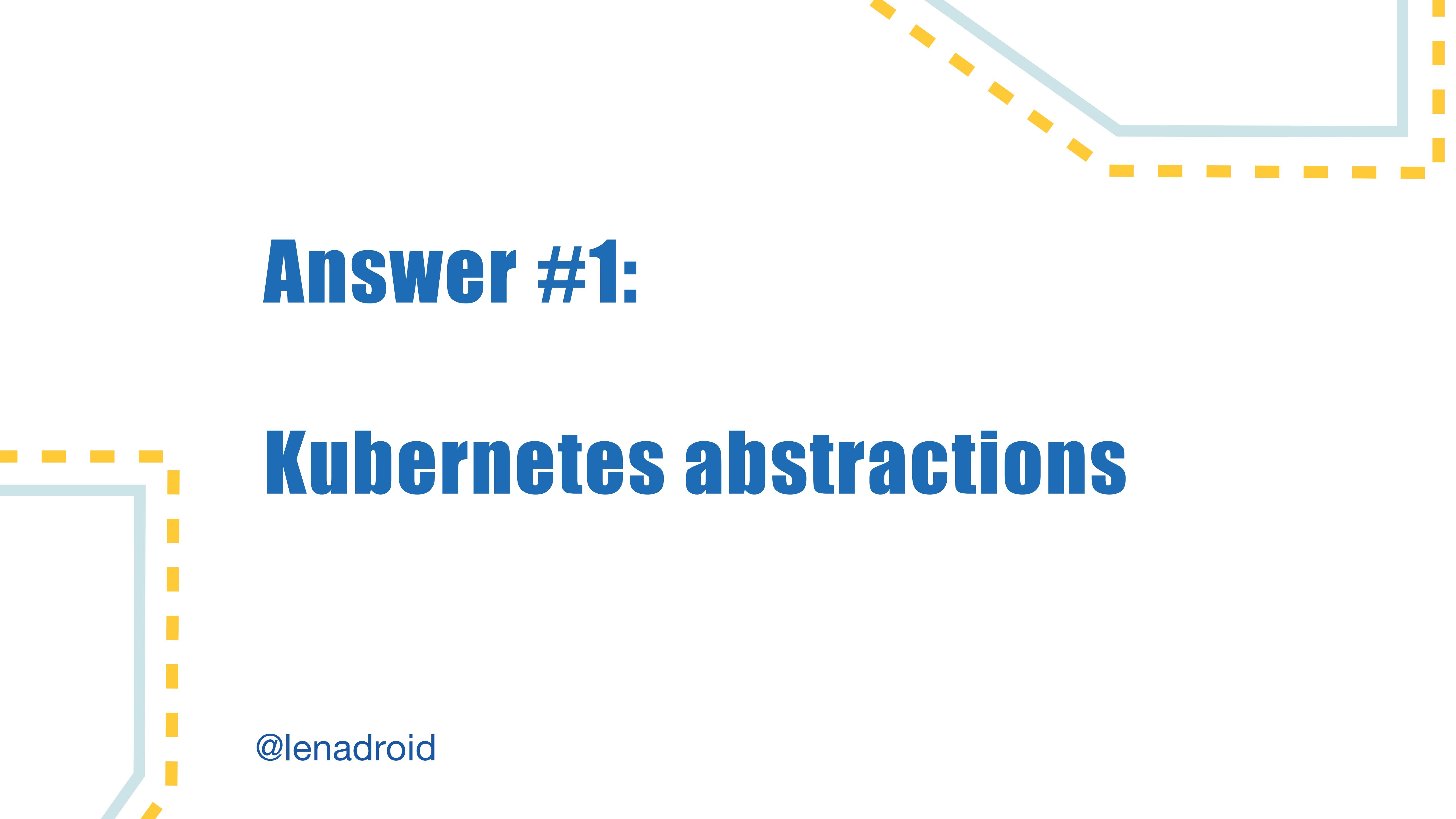
#SeattleScalability

Interpretation #3

**How to get them to run and
work according to what is
considered correct behavior
for those systems?**

@lenadroid

#SeattleScalability

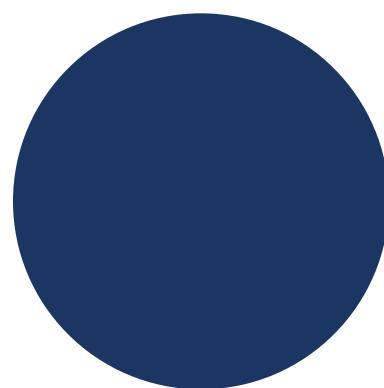


Answer #1:

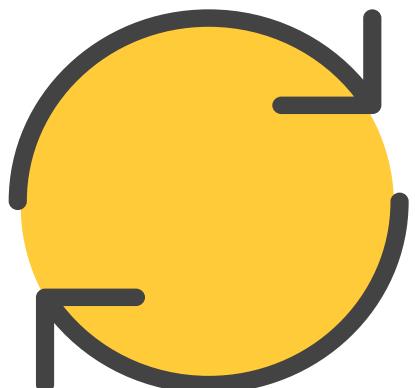
Kubernetes abstractions

@lenandroid

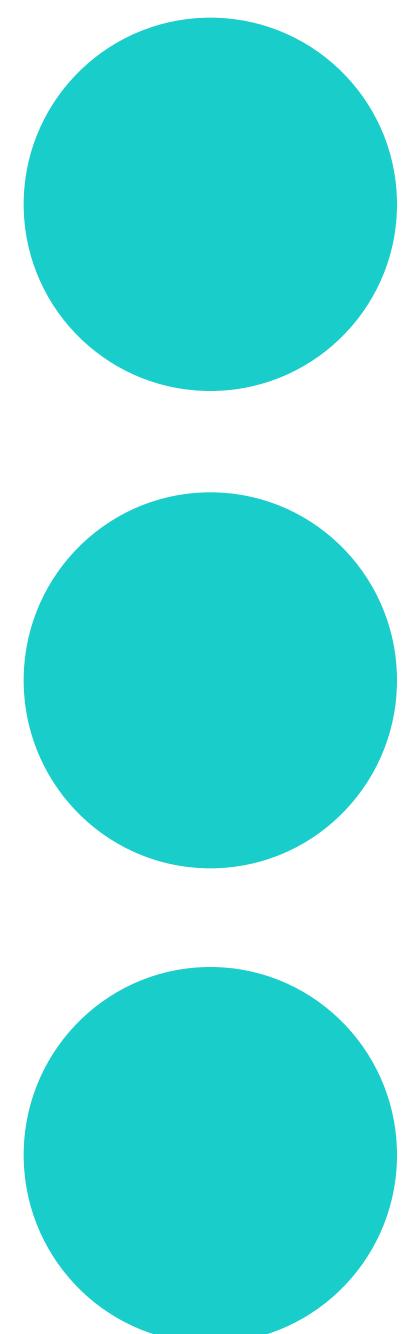
Pod



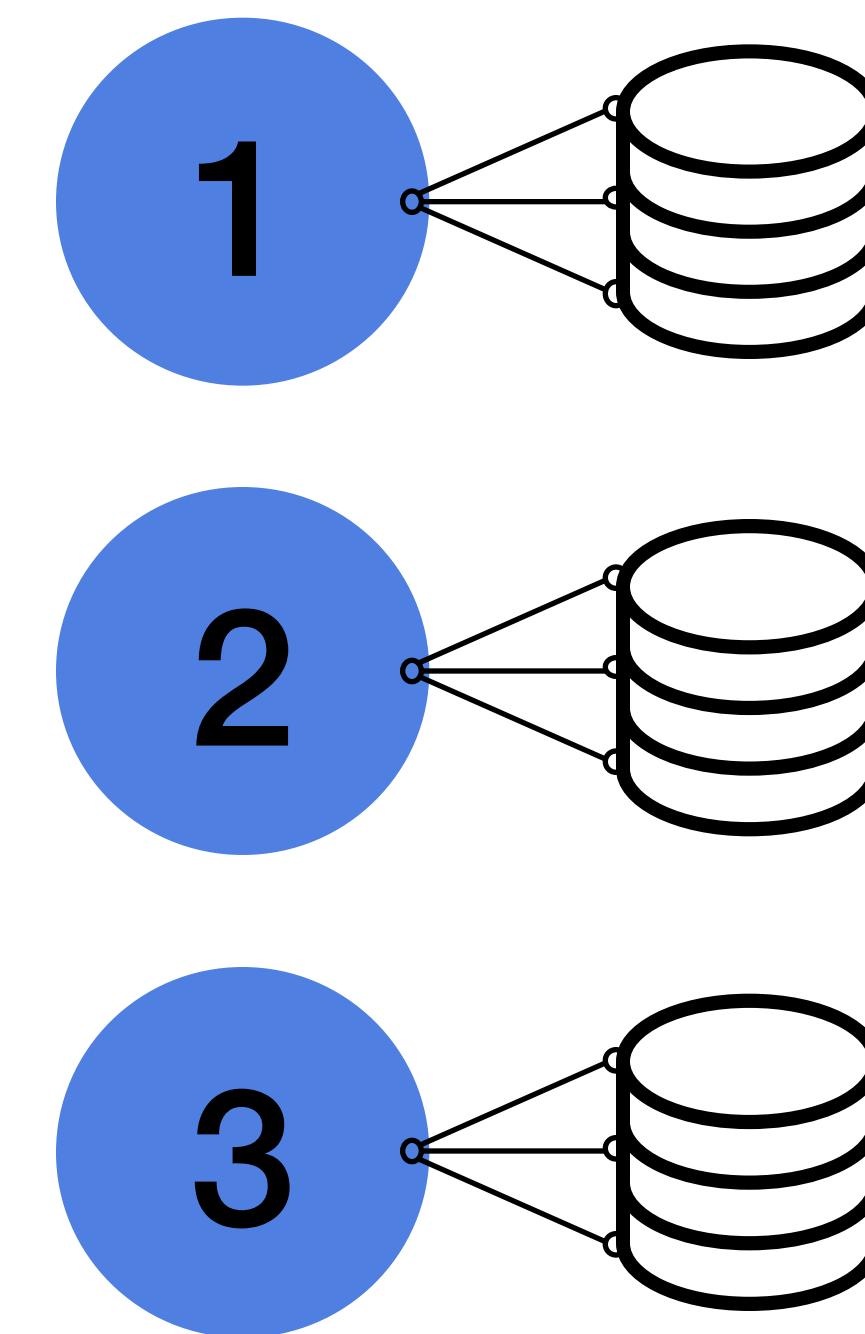
Job



Replica Set

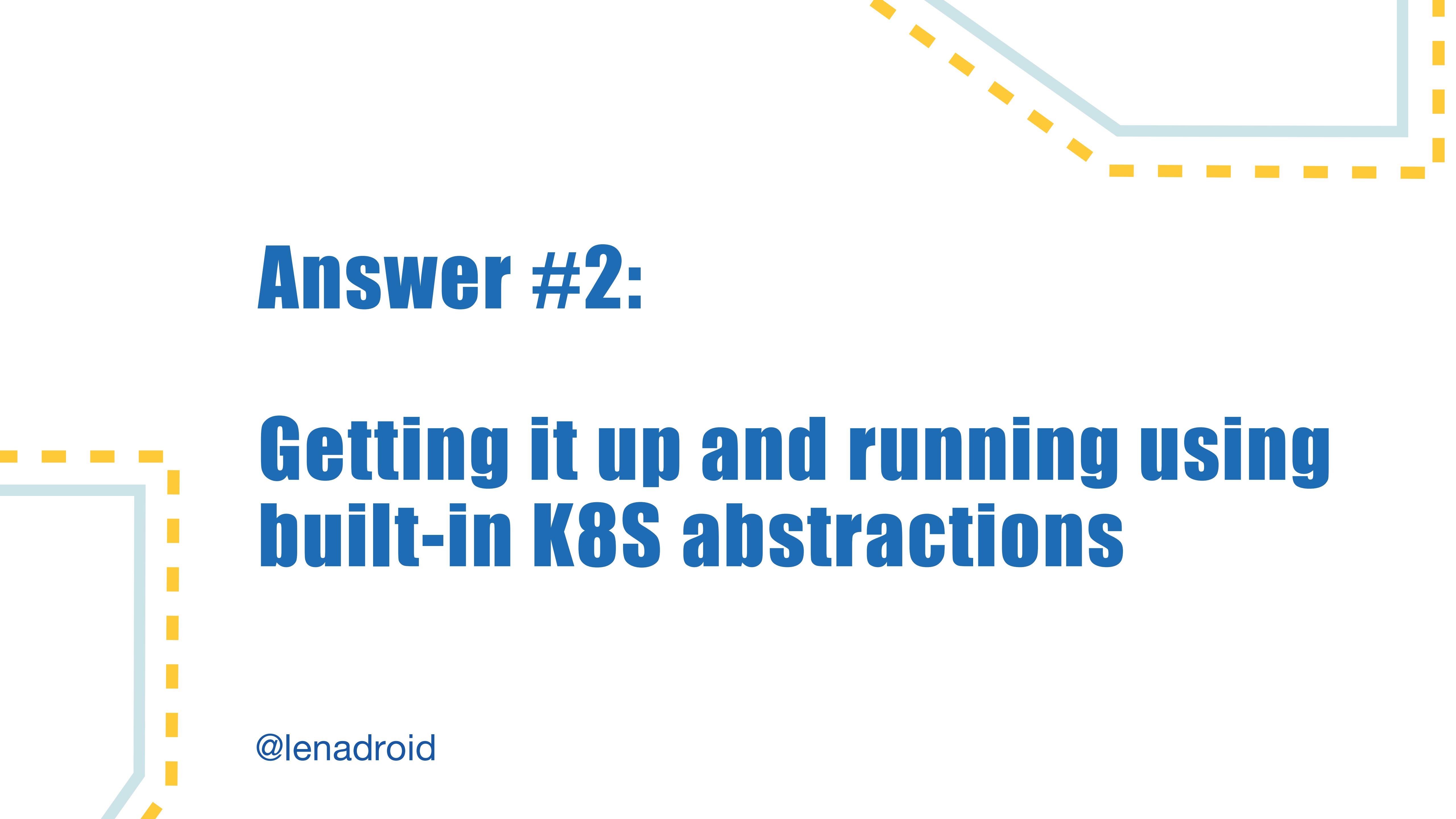


Stateful Set



And many more...

@lenadroid
#SeattleScalability



Answer #2:

Getting it up and running using built-in K8S abstractions

@lenandroid



@lenandroid

Yaml files



@lenandroid

#SeattleScalability

Example:

Getting Cassandra up and running with Stateful Sets and YAMLs on AKS



@lenadroid

#SeattleScalability

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

! cassandra-statefulset.yaml x

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

! cassandra-statefulset.yaml x

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

```
repeat while (true) {
```

UP AND RUNNING

!=

OPERATING CORRECTLY

```
}
```

@lenandroid

#SeattleScalability

[Get Helm](#)[Blog](#)[Docs](#)

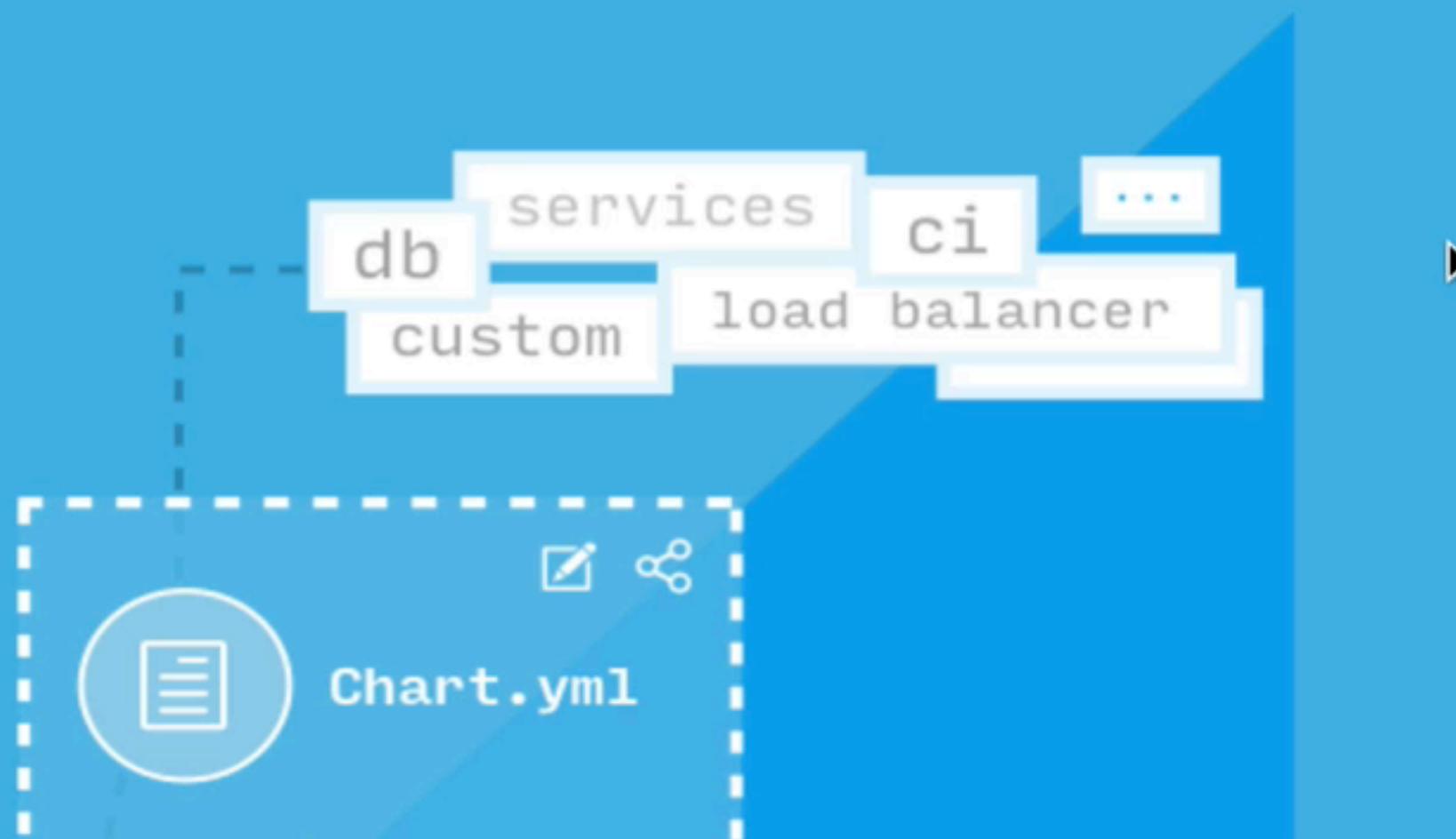
The package manager for Kubernetes

Helm is the best way to find, share, and use software built for [Kubernetes](#).

What is Helm?

Helm helps you manage Kubernetes applications — Helm Charts helps you define, install, and upgrade even the most complex Kubernetes application.

Charts are easy to create, version, share, and publish — so start using Helm and stop the copy-and-paste madness.





Helm

- Charts and Chart.yaml
- Values and Values.yaml
- Chart repositories
- Tiller, Releases and Rollbacks

: @lenadroid
#SeattleScalability



Search or jump to...

Pull requests Issues Marketplace Explore



kubernetes / charts

Watch 164

Star 3,390

Fork 2,899

Code

Issues 397

Pull requests 280

Projects 0

Wiki

Insights

Branch: master

charts / incubator / kafka /

Create new file

Upload files

Find file

History



Dean-Coakley and k8s-ci-robot [incubator/kafka] Fix heapOptions template reference (#5916) ...

Latest commit 3ba5451 a day ago

..



templates

[incubator/kafka] Fix heapOptions template reference (#5916)

a day ago



.helmignore

adds chart for kafka (#144)

a year ago



Chart.yaml

[incubator/kafka] Fix heapOptions template reference (#5916)

a day ago



OWNERS

[incubator/kafka] Allow configurationOverride of zookeeper.connect (#...)

3 months ago



README.md

Add "helm repo add incubator" to README (#5611)

20 days ago



requirements.lock

Update zookeeper dependency to latest chart (#3916)

3 months ago



requirements.yaml

Update zookeeper dependency to latest chart (#3916)

3 months ago



values.yaml

Issue_5426: Add default antiAffinity example to Kafka (#5428)

27 days ago



README.md

Apache Kafka Helm Chart

Chart Details

This chart will do the following:

- Implement a dynamically scalable kafka cluster using Kubernetes StatefulSets
- Implement a dynamically scalable zookeeper cluster as another Kubernetes StatefulSet required for the Kafka cluster above
- Expose Kafka protocol endpoints via NodePort services (optional)

Other examples of deployment

@lenandroid

#SeattleScalability

Example: Spark on Kubernetes (AKS)

@lenandroid

#SeattleScalability

Is this really enough?

@lenadroid
#SeattleScalability

Answer #3:

**Making sure our systems operate
according to what is considered
correct behavior for those systems**

@lenadroid

#SeattleScalability

Things that need special care

@lenadroid
#SeattleScalability

Kubernetes is smart.

P.S. When we tell it how to treat our systems.



A photograph of a modern architectural structure with a highly reflective, curved facade composed of numerous horizontal panels. The building's surface mirrors the surrounding urban environment, including other buildings and possibly a bridge or sky. The lighting suggests it's daytime.

Kafka

Configuration

Rolling cluster restarts and upgrades

Scaling cluster up and down

Data rebalancing

Kafka configuration

Pod specific vs general configurations

Broker ID and rack assignments



Kafka cluster restarts/upgrades

No under-replicated partitions

Restart one broker at a time

Always wait for restarted broker to catch up to a leader

Scaling Kafka cluster

Partition assignment

Data rebalancing

Managing configuration



OPERATORS

@lenandroid

Operators

Custom Resource Definitions

+

Custom Controllers

=

More Control





All

Images

News

Videos

Shopping

More

Settings

Tools

About 69,600 results (0.33 seconds)

GitHub - krallistic/kafka-operator: A Kafka Operator for Kubernetes

<https://github.com/krallistic/kafka-operator> ▾

README.md. kafka-operator - A Kafka Operator for Kubernetes. A Kubernetes Operator for Apache Kafka, which deploys, configures and manages your kafka ...

GitHub - nbogojevic/kafka-operator

<https://github.com/nbogojevic/kafka-operator> ▾

Operator monitors ConfigMap Kubernetes resources that are tagged with config=kafka-topic label. From those ConfigMaps, operator extracts information about ...

Introducing the Confluent Operator: Apache Kafka® on Kubernetes

<https://www.confluent.io/.../introducing-the-confluent-operator-apache-kafka-on-kub...> ▾

May 3, 2018 - With the Confluent Operator, we are productizing years of Apache Kafka experience with Kubernetes expertise to offer our users the best way to ...



All

News

Images

Videos

Shopping

More

Settings

Tools

About 53,600 results (0.39 seconds)

[GitHub - instaclustr/cassandra-operator: Kubernetes operator for ...](#)

<https://github.com/instaclustr/cassandra-operator> ▾

GitHub is where people build software. More than 28 million people use GitHub to discover, fork, and contribute to over 85 million projects.

[GitHub - aslanbekirov/cassandra-operator: cassandra operator ...](#)

<https://github.com/aslanbekirov/cassandra-operator> ▾

The **Cassandra operator** manages **Cassandra** clusters deployed to **Kubernetes** and automates tasks related to operating an **Cassandra** cluster. Create and ...

[GitHub - vgkowski/cassandra-operator: kubernetes operator for ...](#)

<https://github.com/vgkowski/cassandra-operator> ▾

Readme.md. Building the **operator**. **Kubernetes** version: 1.9. This **operator** use the **kubernetes** code-generator for. **clientset**: used to manipulate objects defined ...

Spark [>=2.3]

Automatic job resubmission after specification update

Configurable restart policy

Automatic retries on failed submissions

Mounting specific ConfigMaps and volumes



Example: Spark Operator on Kubernetes (AKS)

@lenandroid

#SeattleScalability

Example: Redis Operator on Kubernetes (AKS)

@lenandroid

#SeattleScalability

Anatomy of an Operator

CRD

Resource Structs/Types

Informer

Controller

Queue

@lenadroid
#SeattleScalability



TensorFlow Operator

@lenadroid
#SeattleScalability

CRD

```
apiVersion: apiextensions.k8s.io/v1beta1
kind: CustomResourceDefinition
metadata:
  name: tfjobs.kubeflow.org
spec:
  group: kubeflow.org
  version: v1alpha1
  names:
    kind: TFJob
    singular: tfjob
    plural: tfjobs
```

TFJOB

```
// TFJob represents the configuration of signal TFJob
type TFJob struct {
    metav1.TypeMeta `json:",inline"`

    // Standard object's metadata.
    metav1.ObjectMeta `json:"metadata,omitempty"`

    // Specification of the desired behavior of the TFJob.
    Spec TFJobSpec `json:"spec,omitempty"`

    // Most recently observed status of the TFJob.
    // This data may not be up to date.
    // Populated by the system.
    // Read-only.
    Status TFJobStatus `json:"status,omitempty"`
}
```

Takeaways

Built-in Kubernetes abstractions don't solve all the issues

Tools like Helm really help structure and manage deployments

Operators = Custom Controllers + CRDs & ...

.... a way to teach Kubernetes understand our needs

There're many ways to deploy and use Operators...

@lenadroid

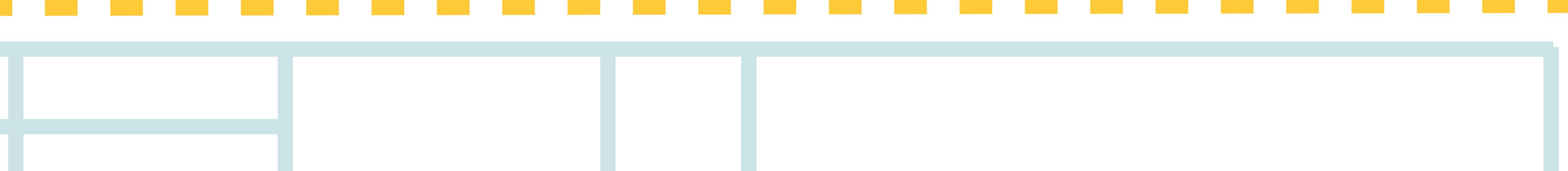
#SeattleScalability

Thank You!

My blog - lenandroid.github.io/posts.html

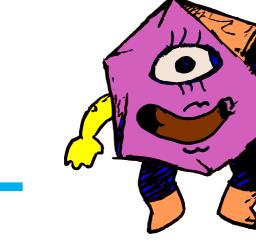
My talk on Stateful Sets, Persistent Volumes, Persistent Volume Claims, Storage Classes - aka.ms/gotochgo

Thomas Stringer's post on Custom Controllers –
aka.ms/custom-controllers



Alena Hall - lenadroid



- ✓ Works on Azure at  Microsoft
- ✓ Lives in  Seattle
- ✓ F# Software Foundation Board of Trustees
- ✓ Organizes [@ML4ALL](#) 
- ✓ Program Committee for Lambda World
- ✓ Has a channel: YouTube /c/AlenaHall