

Distributed systems for big data processing on Kubernetes

Concepts, architecture and implementation

Alena Hall - @lenadroid

Agenda

Distributed systems for big data

Running them on Kubernetes?

Getting them to “work”

Challenges

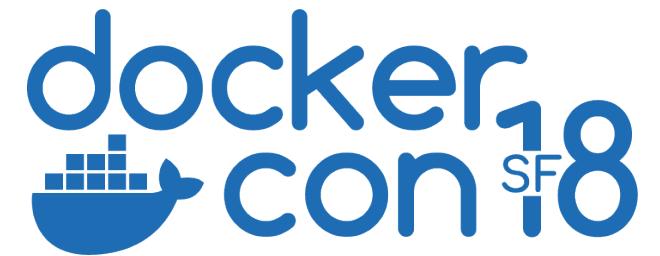
Work in progress and solutions

Future

Summary

Distributed systems for big data: storage, processing, analysis

@lenadroid

 docker
con¹⁸
SF



@lenadroid

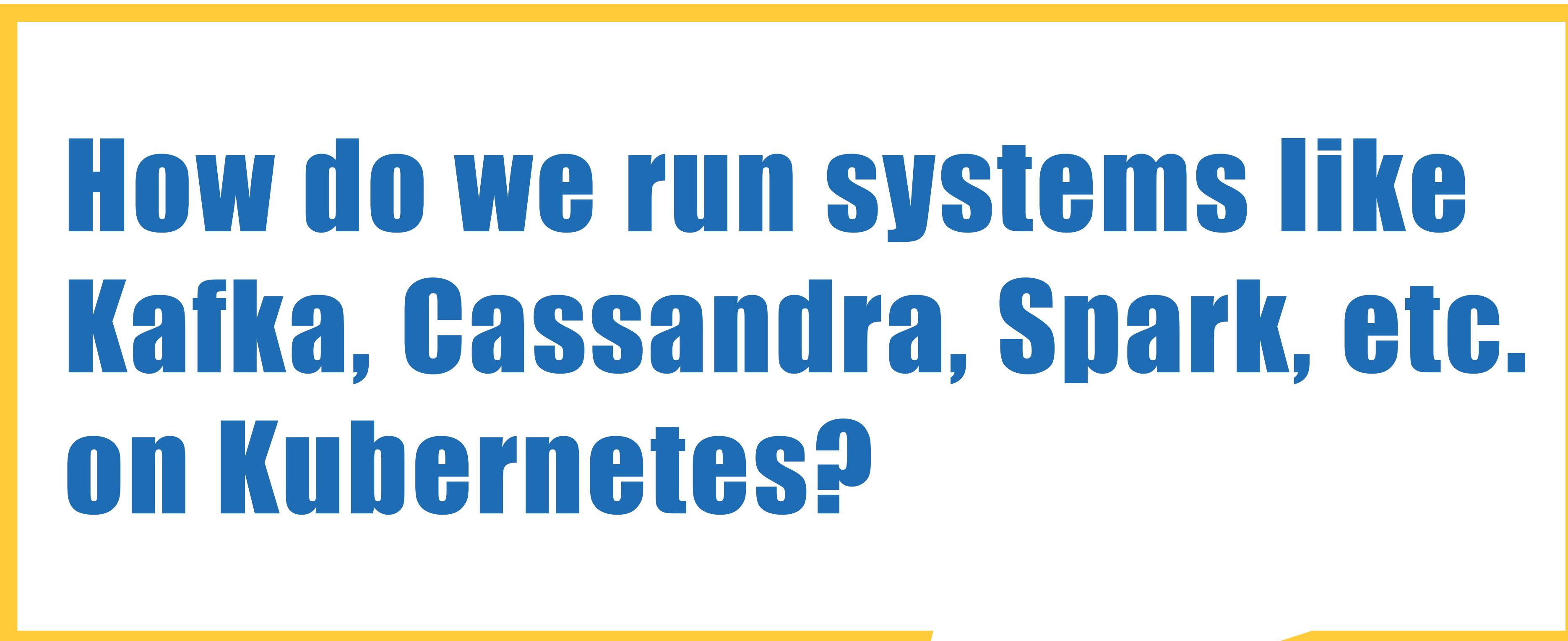


Data systems on Kubernetes



@lenandroid

docker
con¹⁸
SF



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

@lenadroid

 docker
con¹⁸
SF

Interpretation #1

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

@lenadroid

 dockercon¹⁸
SF

Interpretation #2

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

@lenadroid

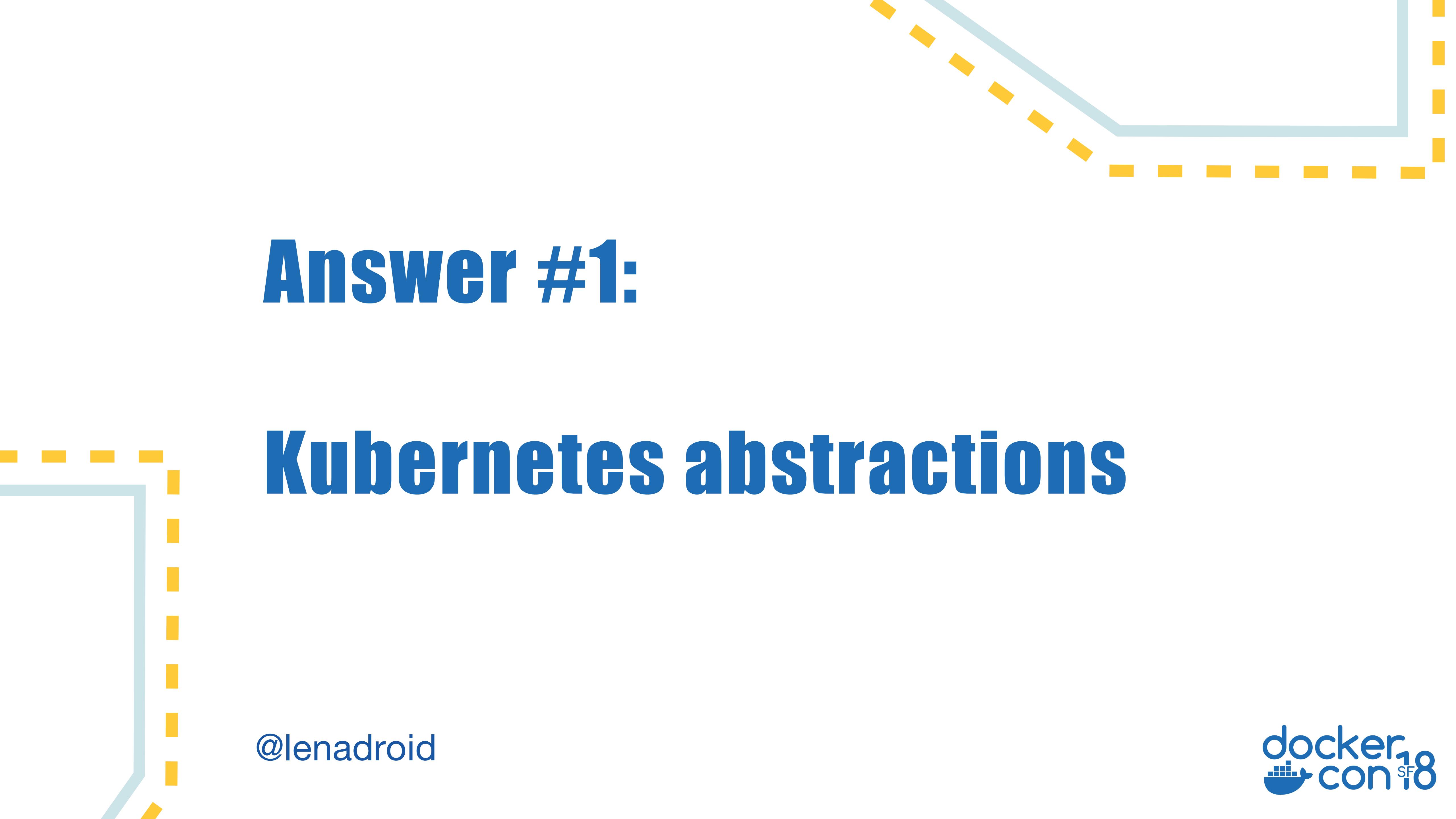
 docker
con¹⁸
SF

Interpretation #3

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

@lenadroid

docker
con¹⁸
SF

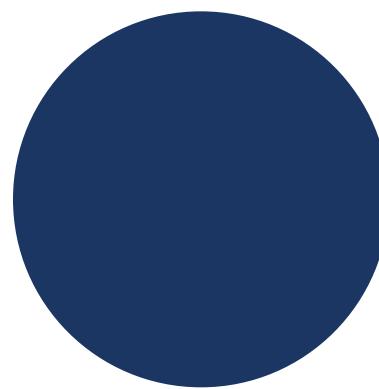


Answer #1:

Kubernetes abstractions

@lenandroid

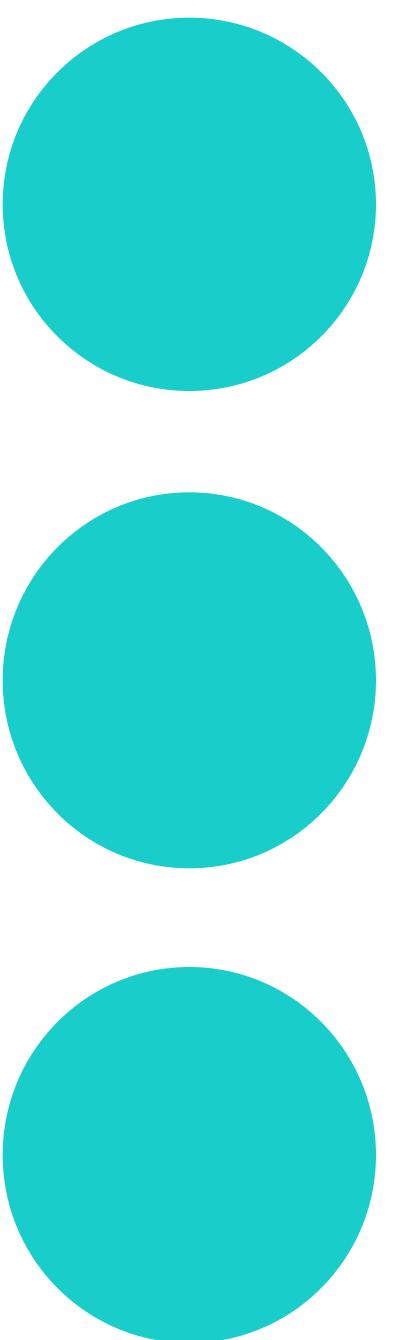
Pod



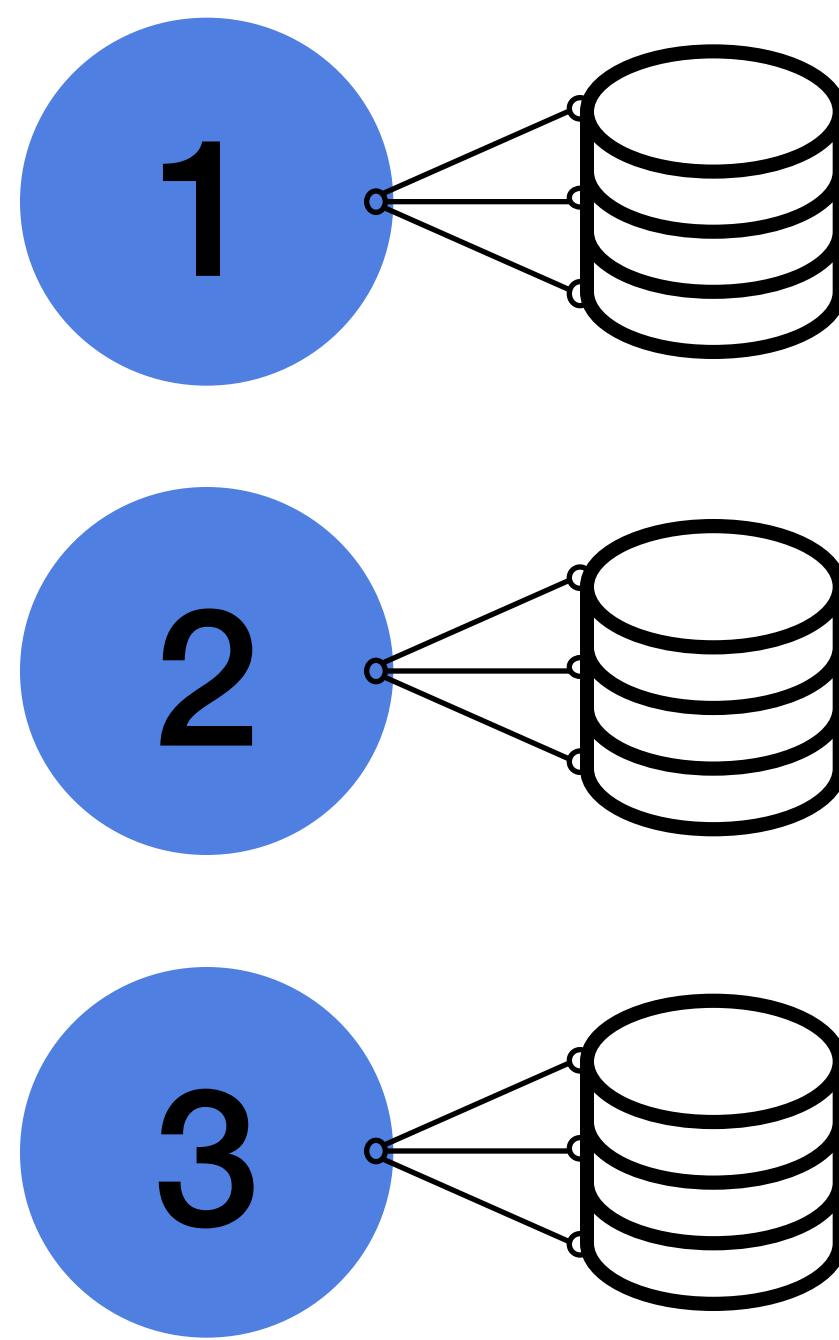
Job



Replica Set



Stateful Set



And many more...

Answer #2:

**Getting it up and running using
built-in K8S abstractions**

@lenandroid

Yaml files



@lenandroid

Example: Stateful Sets + Cassandra

@lenadroid

docker
con SF 18

! cassandra-statefulset.yaml ×

```
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 ×

```
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 ×

```
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

```
}
```

@lenadroid

docker
con
SF 18

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

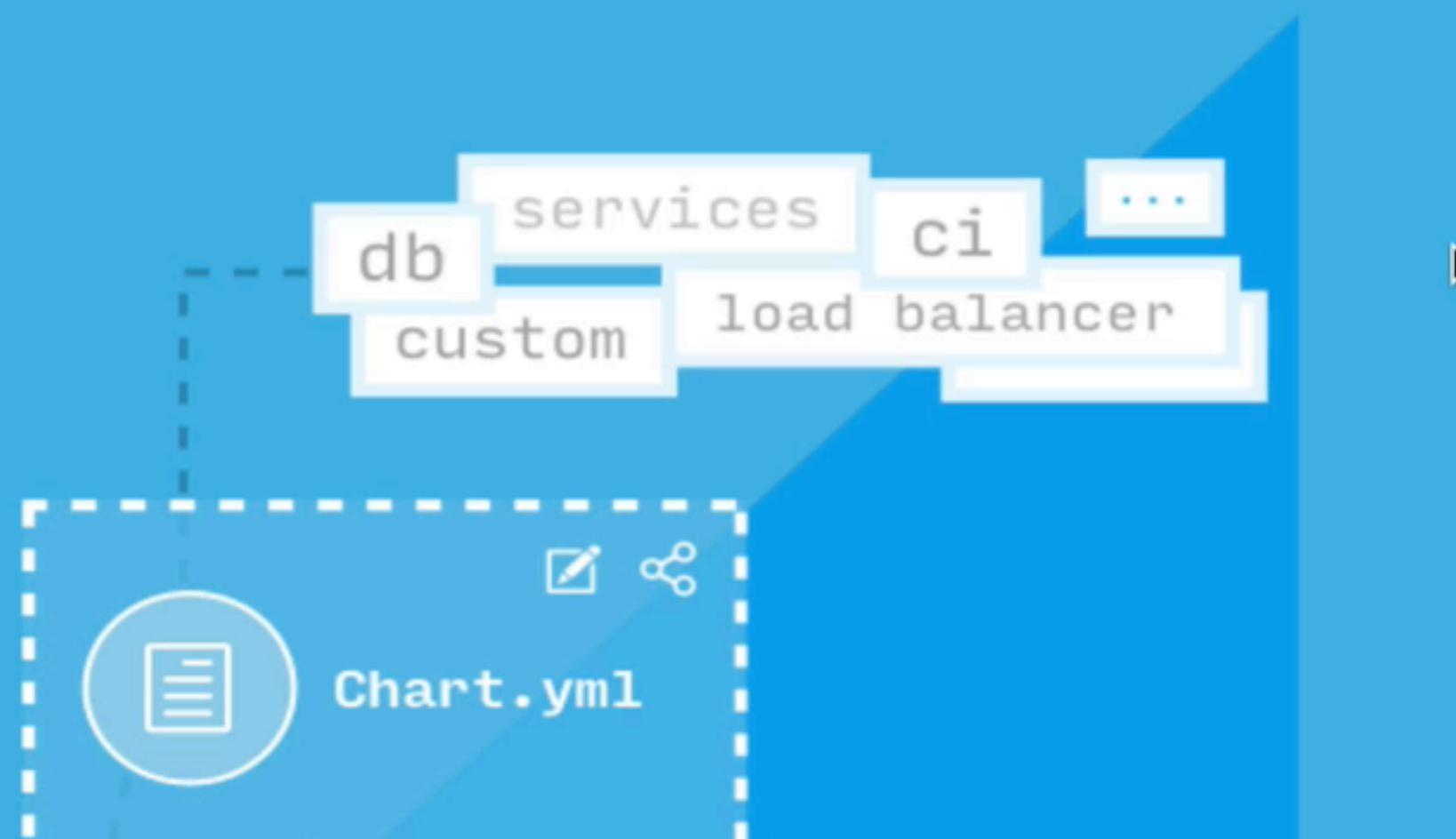
The package manager for Kubernetes

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.



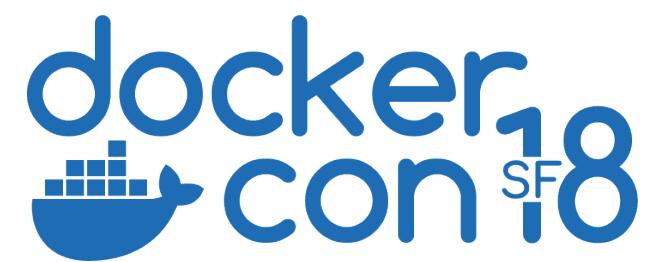
Helm

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



Other examples of deployment

@lenadroid

 docker
con¹⁸
SF

Example: Spark on Kubernetes

@lenadroid

docker
con SF18

Is this really enough?

Answer #3:

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

@lenandroid

docker
con
SF
18

Things that need special care



Kubernetes is smart.

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

Kafka

Configuration

Rolling cluster restarts and upgrades

Scaling cluster up and down

Data rebalancing



Kafka configuration

Pods read configuration from ConfigMaps

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 ConfigMaps



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

@lenadroid

docker
con¹⁸
SF

Anatomy of an Operator

CRD

Resource Structs/Types

Informer

Controller

Queue





TensorFlow Operator

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"`
}
```

```
// TFJobSpec is a desired state description of the TFJob.  
  
type TFJobSpec struct {  
    // TFReplicaSpecs is map of TFReplicaType and TFReplicaSpec  
    // specifies the TF replicas to run.  
    // For example,  
    // {  
    //     "PS": TFReplicaSpec,  
    //     "Worker": TFReplicaSpec,  
    // }  
    TFReplicaSpecs map[TFReplicaType]*TFReplicaSpec `json:"tfReplicaSpecs"  
}
```

```
// TFReplicaSpec is a description of the TFReplica
type TFReplicaSpec struct {
    // Replicas is the desired number of replicas of the given template.
    // If unspecified, defaults to 1.
    Replicas *int32 `json:"replicas,omitempty"`

    // Template is the object that describes the pod that
    // will be created for this TFReplica.
    // We use RestartPolicy in PodTemplateSpec
    // to describe how the containers within the pod should be restarted.
    // Please set this restart policy carefully according to your code.
    Template v1.PodTemplateSpec `json:"template,omitempty"`

    // Restart policy for all TFReplicas within the TFJob.
    // One of Always, OnFailure, Never and ExitCode.
    // Default to Always.
    RestartPolicy RestartPolicy `json:"restartPolicy,omitempty"`
}
```

```
// TFReplicaType is the type for TFReplica.  
type TFReplicaType string  
  
const (  
    // TFReplicaTypePS is the type for parameter servers of distributed TensorFlow.  
    TFReplicaTypePS TFReplicaType = "PS"  
  
    // TFReplicaTypeWorker is the type for workers of distributed TensorFlow.  
    // This is also used for non-distributed TensorFlow.  
    TFReplicaTypeWorker TFReplicaType = "Worker"  
  
    // TFReplicaTypeChief is the type for chief worker of distributed TensorFlow.  
    // If there is "chief" replica type, it's the "chief worker".  
    // Else, worker:0 is the chief worker.  
    TFReplicaTypeChief TFReplicaType = "Chief"  
  
    // TFReplicaTypeEval is the type for evaluation replica in TensorFlow.  
    TFReplicaTypeEval TFReplicaType = "Eval"  
)
```

```
// RestartPolicy describes how the TFReplicas should be restarted.  
// Only one of the following restart policies may be specified.  
// If none of the following policies is specified, the default one  
// is RestartPolicyAlways.  
  
type RestartPolicy string  
  
const (  
    RestartPolicyAlways    RestartPolicy = "Always"  
    RestartPolicyOnFailure RestartPolicy = "OnFailure"  
    RestartPolicyNever     RestartPolicy = "Never"  
  
    // `ExitCode` policy means that user should add exit code by themselves,  
    // `tf-operator` will check these exit codes to  
    // determine the behavior when an error occurs:  
    // - 1-127: permanent error, do not restart.  
    // - 128-255: retryable error, will restart the pod.  
    RestartPolicyExitCode RestartPolicy = "ExitCode"  
)
```

```
// TFJobStatus represents the current observed state of the TFJob.

type TFJobStatus struct {
    // Conditions is an array of current observed TFJob conditions.
    Conditions []TFJobCondition `json:"conditions"`

    // TFReplicaStatuses is map of TFReplicaType and TFReplicaStatus,
    // specifies the status of each TFReplica.
    TFReplicaStatuses map[TFReplicaType]*TFReplicaStatus `json:"tfReplicaStatuses"`

    // Represents time when the TFJob was acknowledged by the TFJob controller.
    // It is not guaranteed to be set in happens-before order across separate operations.
    // It is represented in RFC3339 form and is in UTC.
    StartTime *metav1.Time `json:"startTime,omitempty"`

    // Represents time when the TFJob was completed. It is not guaranteed to
    // be set in happens-before order across separate operations.
    // It is represented in RFC3339 form and is in UTC.
    CompletionTime *metav1.Time `json:"completionTime,omitempty"`

    // Represents last time when the TFJob was reconciled. It is not guaranteed to
    // be set in happens-before order across separate operations.
    // It is represented in RFC3339 form and is in UTC.
    LastReconcileTime *metav1.Time `json:"lastReconcileTime,omitempty"`
}
```

```
type TFJobInformer interface {
    Informer() cache.SharedIndexInformer
    Lister() v1alpha2.TFJobLister
}

type tFJobInformer struct {
    factory internalinterfaces.SharedInformerFactory
}

// NewTFJobInformer constructs a new informer for TFJob type.
// Always prefer using an informer factory to get a shared informer instead of getting an independent
// one. This reduces memory footprint and number of connections to the server.
func NewTFJobInformer(client versioned.Interface, namespace string, resyncPeriod time.Duration, indexers cache.Indexers) cache.SharedIndexI
    return cache.NewSharedIndexInformer(
        &cache.ListWatch{
            ListFunc: func(options v1.ListOptions) (runtime.Object, error) {
                return client.KubeflowV1alpha2().TFJobs(namespace).List(options)
            },
            WatchFunc: func(options v1.ListOptions) (watch.Interface, error) {
                return client.KubeflowV1alpha2().TFJobs(namespace).Watch(options)
            },
        },
        &tensorflow_v1alpha2.TFJob{},
        resyncPeriod,
        indexers,
    )
}
```

```
2 //
3 // Licensed under the Apache License, Version 2.0 (the "License");
4 // you may not use this file except in compliance with the License.
5 // You may obtain a copy of the License at
6 //
7 //     http://www.apache.org/licenses/LICENSE-2.0
8 //
9 // Unless required by applicable law or agreed to in writing, software
10 // distributed under the License is distributed on an "AS IS" BASIS,
11 // WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
12 // See the License for the specific language governing permissions and
13 // limitations under the License.
14
15 // Package controller provides a Kubernetes controller for a TensorFlow job resource.
16 package controller
17
18 import (
19     "errors"
20     "fmt"
21     "time"
22
23     log "github.com/sirupsen/logrus"
24     "k8s.io/api/core/v1"
25     apierrors "k8s.io/apimachinery/pkg/api/errors"
26     metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
27     "k8s.io/apimachinery/pkg/util/runtime"
28     utilruntime "k8s.io/apimachinery/pkg/util/runtime"
29     "k8s.io/apimachinery/pkg/util/wait"
30     "k8s.io/client-go/kubernetes"
31     "k8s.io/client-go/kubernetes/scheme"
32     typedcorev1 "k8s.io/client-go/kubernetes/typed/core/v1"
33     "k8s.io/client-go/tools/cache"
34     "k8s.io/client-go/tools/record"
35     "k8s.io/client-go/util/workqueue"
36
37     "github.com/juju/ratelimit"
38     tfv1alpha1 "github.com/kubeflow/tf-operator/pkg/apis/tensorflow/v1alpha1"
39     tfjobclient "github.com/kubeflow/tf-operator/pkg/client/clientset/versioned"
40     kubeflowscheme "github.com/kubeflow/tf-operator/pkg/client/clientset/versioned/scheme"
```

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



Natallia Dzenisenka

twitter.com/nata_dzen

co-presenter

Thank You!

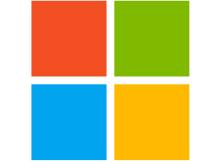
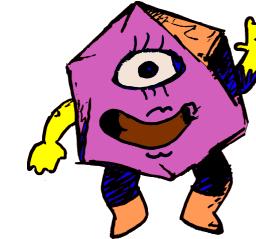
My blog - lenandroid.github.io

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](#)