

# Knative Tutorial

Red Hat Developer Experience Team

Version 0.0.1

# Table of Contents

Overview	1
Setup	2
Prerequisite CLI tools	2
Download Tutorial Sources	2
Minishift	3
Enable Admission Controller Webhook	3
Knative Setup	4
Configuring OpenShift project for Knative applications	9
Work folder	9
Basics and Fundamentals	10
Prerequisite	10
Build Containers	11
Deploy Service	11
Invoke Service	13
See what you have deployed	13
Deploy new Revision Service	14
Pinning service to a revision	16
Cleanup	17
Configurations and Routes	18
Prerequisite	18
Build Containers	18
Deploy Configuration	19
Invoke Service	20
Deploy Route	21
See what you have deployed	22
Deploy new Revision Service	23
Distributing traffic	24
Cleanup	26
Scaling	27
Prerequisite	27
Build Containers	28
Deploy Service	28
Invoke Service	30
Auto Scaling	33
Minimum Scale	35
Cleanup	37
Build	38
Prerequisite	38

Generate Docker Secret .....	39
Generate Knative build spec .....	40
Generate Knative service .....	42
Create build .....	43
Watching logs .....	43
Checking build status .....	44
Deploy service using Build .....	45
Invoke service .....	46
Cleanup .....	46
Build Template .....	47
Prerequisite .....	47
Generate build template .....	48
Generate Knative service .....	50
Apply resources .....	51
Create build template .....	52
See what you have deployed .....	52
Invoke Service .....	53
Cleanup .....	53
FAQs .....	55
How to access the Knative services ? .....	55
Why dev.local suffixes for container images? .....	55
What is revision simpler terms? .....	57

# Overview

Serverless epitomize the very benefits of what cloud platforms promise: offload the management of infrastructure while taking advantage of a consumption model for the actual utilization of services. While there are a number of server frameworks out there, Knative is the first serverless platform specifically designed for kubernetes and OpenShift.

Serverless computing refers to the concept of building and running applications that do not require server management. It describes a finer-grained deployment model where applications, bundled as one or more functions, are uploaded to a platform and then executed, scaled, and billed in response to the exact demand needed at the moment”

**Source:** <https://www.cncf.io/blog/2018/02/14/cncf-takes-first-step-towards-serverless-computing/>

This tutorial will act as step-by-step guide in helping you to understand Knative starting with setup, understanding fundamentals concepts such as service, configuration, revision etc., and finally deploying some use cases which could help deploying serverless applications at enterprises.

This content is brought to you by [Red Hat Developer Program](#) - Register today!

# Setup

## Prerequisite CLI tools

You will need in this tutorial:

Tool	macOS	Fedora
<code>minishift</code>	<a href="https://github.com/minishift/minishift/releases">https://github.com/minishift/minishift/releases</a>	<a href="https://github.com/minishift/minishift/releases">https://github.com/minishift/minishift/releases</a>
<code>docker</code>	<a href="#">Docker for Mac</a>	<code>dnf install docker</code>
<code>kubectl</code>	<a href="#">Mac OS</a>	<code>dnf install kubernetes-client</code>
<a href="#">Apache Maven</a>	<code>brew install maven</code>	<code>dnf install maven</code>
<code>stern</code>	<code>brew install stern</code>	<code>sudo curl --output /usr/local/bin/stern -L https://github.com/wercker/stern/releases/download/1.6.0/stern_linux_amd64 &amp;&amp; sudo chmod +x /usr/local/bin/stern</code>
<code>curl</code> , <code>gunzip</code> , <code>tar</code>	built-in or part of your bash shell	should also be installed already, but just in case... <code>dnf install curl gzip tar</code>
<code>git</code>	<code>brew install git</code>	<code>dnf install git</code>
<code>xmlstarlet</code>	<code>brew install xmlstarlet</code>	<a href="http://xmlstar.sourceforge.net/">http://xmlstar.sourceforge.net/</a>
<code>yq</code>	<code>brew install yq</code>	<a href="https://github.com/mikefarah/yq/releases/latest">https://github.com/mikefarah/yq/releases/latest</a>
<code>httpie</code>	<code>brew install httpie</code>	<code>dnf install httpie</code>
<code>go</code>	<code>brew install golang</code>	<a href="https://dl.google.com/go/go1.11.5.linux-amd64.tar.gz">https://dl.google.com/go/go1.11.5.linux-amd64.tar.gz</a>
<code>hey</code>	<code>go get -u github.com/rakyll/hey</code>	<code>go get -u github.com/rakyll/hey</code>
<a href="https://jsonnet.org/">https://jsonnet.org/</a>	<code>brew install jsonnet</code>	<a href="https://github.com/google/jsonnet">https://github.com/google/jsonnet</a>

## Download Tutorial Sources

Before we start to setting up the environment, lets clone the tutorial sources and call the cloned folder as `$TUTORIAL_HOME`

```
git clone https://github.com/kameshsampath/knative-tutorial
```

The `work` folder in `$TUTORIAL_HOME` can be used to download the demo application resources and refer them during the exercises. The `work` folder has a README, which has instructions on source code repo and git commands to clone the sources.



All examples in this tutorial are tested with **OpenShift 3.11** and **Knative 0.3.0**

## Minishift

All the demos in this tutorial will be run using [minishift](#), the single node [OKD](#)(the Origin Community Distribution of kubernetes that powers Red Hat OpenShift) cluster that can be used for development purposes.

### Setup minishift

```
export TIMEZONE=$(sudo systemsetup -gettimezone | awk '{print $3}') ①
minishift profile set knative-tutorial
minishift config set openshift-version v3.11.0
minishift config set memory 8GB # if you more then preferred is 10GB
minishift config set cpus 4
minishift config set disk-size 50g
minishift config set image-caching true
minishift addons enable admin-user

minishift start --timezone $TIMEZONE

eval $(minishift docker-env) && eval $(minishift oc-env)
```

① Some examples we use in this tutorial are time sensitive, hence its recommended to set the appropriate timezone. If you are linux please use the appropriate command to get the timezone of the host.

## Enable Admission Controller Webhook

```
#!/bin/bash

minishift openshift config set --target=kube --patch '{
  "admissionConfig": {
    "pluginConfig": {
      "ValidatingAdmissionWebhook": {
        "configuration": {
          "apiVersion": "apiserver.config.k8s.io/v1alpha1",
          "kind": "WebhookAdmission",
          "kubeConfigFile": "/dev/null"
        }
      },
      "MutatingAdmissionWebhook": {
        "configuration": {
          "apiVersion": "apiserver.config.k8s.io/v1alpha1",
          "kind": "WebhookAdmission",
          "kubeConfigFile": "/dev/null"
        }
      }
    }
  }
}'

# Allowing few minutes for the OpenShift to be restarted
until oc login -u admin -p admin 2>/dev/null; do sleep 5; done;
```

## Knative Setup

### Login as admin user

Since there was an `admin` addon added during minishift configuration, its now possible to login with the `admin` user. For example:

```
oc login -u admin -p admin
```

### Install Istio

#### Configure Istio Policies

Knative depends on Istio, to deploy Istio you need to run the following commands to configure necessary [privileges](#) to the service accounts used by Istio.

```
oc adm policy add-scc-to-user anyuid -z istio-ingress-service-account -n istio-system
oc adm policy add-scc-to-user anyuid -z default -n istio-system
oc adm policy add-scc-to-user anyuid -z prometheus -n istio-system
oc adm policy add-scc-to-user anyuid -z istio-egressgateway-service-account -n istio-system
oc adm policy add-scc-to-user anyuid -z istio-citadel-service-account -n istio-system
oc adm policy add-scc-to-user anyuid -z istio-ingressgateway-service-account -n istio-system
oc adm policy add-scc-to-user anyuid -z istio-cleanup-old-ca-service-account -n istio-system
oc adm policy add-scc-to-user anyuid -z istio-mixer-post-install-account -n istio-system
oc adm policy add-scc-to-user anyuid -z istio-mixer-service-account -n istio-system
oc adm policy add-scc-to-user anyuid -z istio-pilot-service-account -n istio-system
oc adm policy add-scc-to-user anyuid -z istio-sidecar-injector-service-account -n istio-system
oc adm policy add-scc-to-user anyuid -z cluster-local-gateway-service-account -n istio-system
oc adm policy add-scc-to-user anyuid -z istio-galley-service-account -n istio-system
```

## Deploy Istio Components

```
curl -L https://github.com/knative/serving/releases/download/v0.3.0/istio.yaml \
| sed 's/LoadBalancer/NodePort/' \
| kubectl apply --filename -
```



The above command throws some warnings and error as to some objects not found, it can safely ignore or run the command again

You can monitor the Istio components via the command:

```
oc get pods -n istio-system -w
```



It will take a few minutes for all the Istio components to be up and running. Please wait for all the Istio pods to be running before deploying [Knative Serving](#). Use **CTRL+C** to exit watch mode.

A successful deployment should show similar output as shown below.



```
$ oc get pods -n istio-system
```

NAME	READY	STATUS	RESTARTS	AGE
cluster-local-gateway-76db55c785-9f7zc	1/1	Running	0	2m
istio-citadel-746c765786-fqx8m	1/1	Running	0	2m
istio-cleanup-secrets-znzc6	0/1	Completed	0	2m
istio-egressgateway-7b46794587-m6tqm	1/1	Running	0	2m
istio-galley-75c6976d79-bcnbp	1/1	Running	0	2m
istio-ingressgateway-57f76dc4db-85cps	1/1	Running	0	2m
istio-pilot-6495978c49-6jrmf	2/2	Running	0	2m
istio-pilot-6495978c49-7l6w7	2/2	Running	0	2m
istio-pilot-6495978c49-fh9n8	2/2	Running	0	2m
istio-policy-6677c87b9f-7wxj	2/2	Running	0	2m
istio-sidecar-injector-879fd9dfc-82wjl	1/1	Running	0	2m
istio-statsd-prom-bridge-549d687fd9-mlnb2	1/1	Running	0	2m
istio-telemetry-7d46d668db-4q wz5	2/2	Running	0	2m

## Update Istio sidecar injector ConfigMap

The Istio v1.0.1 release automatic sidecar injection has removed `privileged:true` from init containers, this will cause the Pods with istio proxies automatic inject to crash. Run the following command to update the **istio-sidecar-injector** ConfigMap.

The following command ensures that the `privileged:true` is added to the **istio-sidecar-injector** ConfigMap:

```
oc apply -n istio-system -f $TUTORIAL_HOME/patches/istio-sidecar-injector.yaml
```



Run the above command only once per minishift instance

## Install Knative Build

```
# Setup Knative Build Policies
oc adm policy add-scc-to-user anyuid -z build-controller -n knative-build
oc adm policy add-scc-to-user anyuid -z build-controller -n knative-build

# Install Knative Build components
kubectl apply --filename
https://github.com/knative/build/releases/download/v0.3.0/release.yaml

# give cluster admin privileges to Service Account Build Controller on project
knative-build
oc adm policy add-cluster-role-to-user cluster-admin -z build-controller -n knative-build
oc adm policy add-cluster-role-to-user cluster-admin -z build-controller -n knative-build
```

```
oc get pods -n knative-build -w
```



It will take a few minutes for all the Knative Build components to be up and running. Use **CTRL+C** to exit watch mode.

A successful deployment should show similar output as shown below.

```
$ oc get pods -n knative-build
```

NAME	READY	STATUS	RESTARTS	AGE
build-controller-79cb969d89-bnzbr	1/1	Running	0	28s
build-webhook-58d685fc58-f5s4l	1/1	Running	0	27s

## Install Knative Serving

```
# Setup Knative Serving Policies
oc adm policy add-scc-to-user anyuid -z controller -n knative-serving
oc adm policy add-scc-to-user anyuid -z autoscaler -n knative-serving

# Install Knative Serving components
curl -L https://github.com/knative/serving/releases/download/v0.3.0/serving.yaml \
| sed 's/LoadBalancer/NodePort/' \
| kubectl apply --filename -

# give cluster admin privileges to Service Account Controller on project knative-
serving
oc adm policy add-cluster-role-to-user cluster-admin -z controller -n knative-serving
```

You can monitor the Knative Serving components via components via the command:

```
oc get pods -n knative-serving -w
```



It will take a few minutes for all the Knative Serving components to be up and running. Use **CTRL+C** to exit watch mode.

A successful deployment should show similar output as shown below.

```
$ oc get pods -n knative-serving
```

NAME	READY	STATUS	RESTARTS	AGE
activator-598b4b7787-sb9gv	2/2	Running	0	35s
autoscaler-5cf5cfb4dc-9zd7p	2/2	Running	0	34s
controller-7fc84c6584-ggp7b	1/1	Running	0	38s
webhook-7797ffb6bf-swr4c	1/1	Running	0	38s

## Install Knative Eventing

```
# Setup Knative Eventing Policies
oc adm policy add-scc-to-user anyuid -z eventing-controller -n knative-eventing
oc adm policy add-scc-to-user anyuid -z in-memory-channel-dispatcher -n knative-eventing
oc adm policy add-scc-to-user anyuid -z in-memory-channel-controller -n knative-eventing

# Install Knative Eventing components
kubectl apply --filename
https://github.com/knative/eventing/releases/download/v0.3.0/release.yaml
kubectl apply --filename https://github.com/knative/eventing-
sources/releases/download/v0.3.0/release.yaml

# give cluster admin privileges to Service Accounts on project knative-eventing
oc adm policy add-cluster-role-to-user cluster-admin -z eventing-controller -n
knative-eventing
oc adm policy add-cluster-role-to-user cluster-admin -z default -n knative-sources
oc adm policy add-cluster-role-to-user cluster-admin -z in-memory-channel-dispatcher
-n knative-eventing
oc adm policy add-cluster-role-to-user cluster-admin -z in-memory-channel-controller
-n knative-eventing
```

You can monitor the Knative Eventing components via components via the command:

```
oc get pods -n knative-eventing -w
oc get pods -n knative-sources -w
```



It will take a few minutes for all the Knative Eventing components to be up and running. Use **CTRL+C** to exit watch mode.

```
$ oc get pods -n knative-eventing
NAME                                READY    STATUS    RESTARTS   AGE
eventing-controller-847d8cf969-zv77v 1/1      Running   0           1m
in-memory-channel-controller-59dd7cfb5b-s82x6 1/1      Running   0           1m
webhook-7cfff8d86d-4kljw             1/1      Running   0           1m

$ oc get pods -n knative-sources
NAME                READY    STATUS    RESTARTS   AGE
controller-manager-0 1/1      Running   0           1m
```

# Configuring OpenShift project for Knative applications

```
oc new-project knativetutorial
oc adm policy add-scc-to-user privileged -z default ①
oc adm policy add-scc-to-user anyuid -z default
```

① The `oc adm policy` adds the **privileged** [Security Context Constraints\(SCCs\)](#) to the **default** Service Account. The SCCs are the precursor to the PSP (Pod Security Policy) mechanism in kubernetes.

(OR)

```
kubectl create namespace knativetutorial
```

## Work folder

The work folder i.e `$TUTORIAL_HOME/work` can be used as a work directory during the build. The README in the work folder as the GitHub repository links of the applications `greeter` and `event-greeter` that will be used in various exercises.

# Basics and Fundamentals

At the end of this chapter you will be able to understand and know how to :

- Deploy first Knative service ?
- Deploying multiple revisions of the service
- What different service deployment strategies possible (service vs configurations/routes) ?
- Always running latest vs pinning revision to services

## Prerequisite

The following checks ensures that each chapter exercises are done with right environment settings.

- Make sure to be connected to minishift docker daemon and using right openshift client.

```
eval $(minishift docker-env) && eval $(minishift oc-env)
```

- Make sure right OpenShift and kubernetes client are used.

```
# oc v3.11.0+0cbc58b kubernetes v1.11.0+d4cacc0
oc version
```

(OR)

```
# kubernetes v1.11+
kubectl version
```

- Make sure right maven version is used (only for Java exercises)

```
# right maven version Apache Maven 3.5.x or above
./mvnw --version
```

- Make sure to be on **knativetutorial** OpenShift project/kubernetes namespace

```
# verify to return knativetutorial
oc project -q
```

If you are not on **knativetutorial** project, use the command **oc project knativetutorial** to change to **knativetutorial** project

# Build Containers

- Navigate to work folder

```
cd $TUTORIAL_HOME/work
```

- Clone the application source code

```
git clone https://github.com/redhat-developer-demos/knative-tutorial-greeter.git
```

- Build the application and container image

```
cd knative-tutorial-greeter/java/springboot
./mvnw clean package
docker build -t dev.local/rhdevelopers/greeter:0.0.1 .
```

- Check if images are built and available:

```
docker images | grep dev.local
```

The above command should return something like as shown below:

```
# dev.local/rhdevelopers/greeter 0.0.1 6aa8771da8dd 2 hours ago
456MB
```



You can also use the script `buildImage.sh` that will build source and create the container image

## Deploy Service

Navigate to folder `$TUTORIAL_HOME/01-basics/knative`.

The following snippet shows how a Knative service YAML will look like:

```

apiVersion: serving.knative.dev/v1alpha1
kind: Service
metadata:
  name: greeter
spec:
  runLatest: ①
  configuration:
    revisionTemplate:
      spec:
        container:
          image: dev.local/rhdevelopers/greeter:0.0.1 ②

```

### service.yaml

- ① Makes Knative to always run the latest revision of the deployment
- ② It is very important that the image is a fully qualified name docker image name with tag. For more details on this [Question 2 of FAQ](#)

The service could be deployed using the command:

```
oc apply -n knativetutorial -f service.yaml
```

(OR)

```
oc apply -n knativetutorial -f service.yaml
```

After successful deployment of the service we should see a kubernetes deployment called **greeter-00001-deployment** available.

The screenshot shows the OpenShift Kubernetes Dashboard interface. On the left is a sidebar with navigation options: Overview, Applications, Builds, Resources, Storage, Monitoring, and Catalog. The main panel displays the 'greeter-00001' application. Under the 'DEPLOYMENT' section, it shows 'greeter-00001-deployment, #1'. The 'CONTAINERS' section lists two containers: 'user-container' with image 'dev.local/greeter:1.0.0' and port '8080/TCP (user-port)', and 'queue-proxy' with image 'dev.local/greeter:1.0.0' and port '8012/TCP (queue-port) and 1 other'. The 'NETWORKING' section shows 'Service - Internal Traffic' for 'greeter-00001-service' and '80/TCP (http) → queue-port'. On the right, a circular indicator shows '1 pod'.

# Invoke Service

```
INGRESSGATEWAY=istio-ingressgateway  
IP_ADDRESS="$(minishift ip):$(kubectl get svc $INGRESSGATEWAY --namespace istio-system  
--output 'jsonpath={.spec.ports[?(@.port==80)].nodePort}')" 
```

```
curl -H "Host: greeter.knativetutorial.example.com" $IP_ADDRESS
```

(OR)

```
http $IP_ADDRESS 'Host: greeter.knativetutorial.example.com'
```

The last curl command should return a response like **Hi greeter** ⇒ **greeter-00001-deployment-5d696cc6c8-m65s5: 1**



Sometimes the response might not be returned immediately especially when the pod is coming up from dormant state, at those times try giving request again

## See what you have deployed

The service based deployment strategy that we did now will create many knative resources, the following commands will help you to query and find what has been deployed.

### service

```
# get a Knative Service (short name ksvc) called greeter  
oc -n knativetutorial get services.serving.knative.dev greeter
```

(OR)

```
kubectl -n knativetutorial get services.serving.knative.dev greeter
```

### configuration

```
# get a Knative configuration called greeter  
oc -n knativetutorial get configurations.serving.knative.dev greeter
```

(OR)

```
kubectl -n knativetutorial get configurations.serving.knative.dev greeter
```



## routes

```
# get a Knative routes called greeter
oc -n knativetutorial get routes.serving.knative.dev greeter
```

(OR)

```
kubectl -n knativetutorial get routes.serving.knative.dev greeter
```

When the service was invoked with `curl -H "Host: greeter.knativetutorial.example.com" $IP_ADDRESS`, you noticed that we added a **Host** header to the request with value `greeter.knativetutorial.example.com`, this FQDN is automatically assigned to your Knative service by the Knative Routes, it uses the format like `<service-name>.<namespace>.<domain-suffix>`.



- The domain suffix in this case *example.com* is configurable via the config map **config-domain** of **knative-serving** namespace.

## revisions

```
# get a Knative revisions called greeter you will see only one like greeter-00001
oc -n knativetutorial get revisions.serving.knative.dev
```

(OR)

```
kubectl -n knativetutorial get revisions.serving.knative.dev
```



- add `-oyaml` to the commands above to see more details

## Deploy new Revision Service

As Knative follows [12-Factor](#) application principles, any new [configuration](#) change will trigger new revision of the deployment.

To deploy a new revision of the greeter service, we will add an environment variable to the existing service as shown below:

### Service revision 2

```

apiVersion: serving.knative.dev/v1alpha1
kind: Service
metadata:
  name: greeter
spec:
  runLatest:
    configuration:
      revisionTemplate:
        spec:
          container:
            image: dev.local/rhdevelopers/greeter:0.0.1
            env:
              - name: MESSAGE_PREFIX ①
                value: Hello

```

### service-env.yaml

① Adding an environment variable that will be used as a message prefix

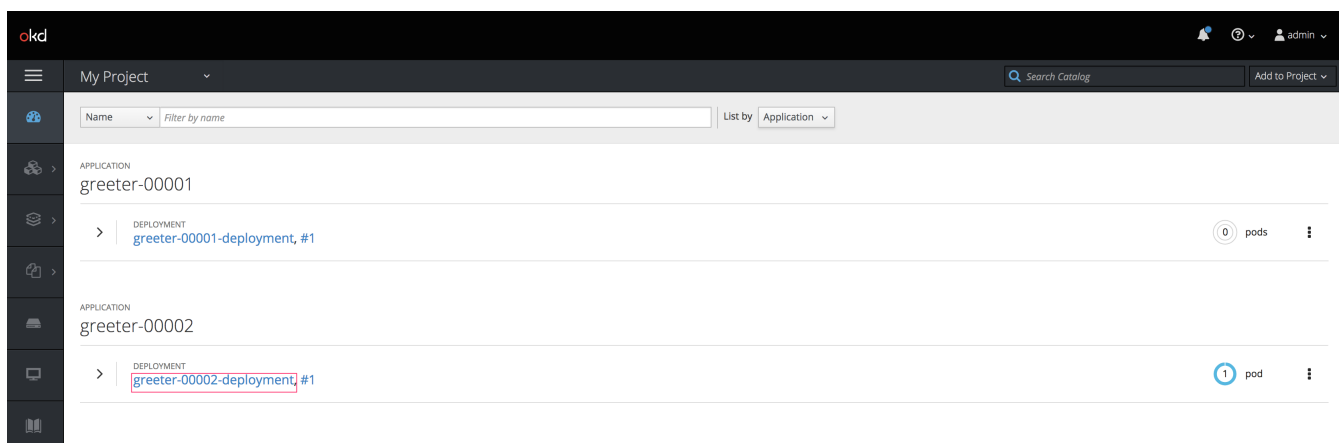
Let us deploy the new revision using the command:

```
oc apply -n knativetutorial -f service-env.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f service-env.yaml
```

After successful deployment of the service we should see a kubernetes deployment called **greeter-00002-deployment** available in the OpenShift dashboard:



Now running the **command** will show two revisions namely **greeter-00001** and **greeter-00002**.

**Invoking Service** will now show an output like **Hello greeter ⇒ greeter-00002-deployment-8d9984dc8-rgzx6: 2**, where **Hello** is the value that we configured via environment variable in the Knative service resource file.

# Pinning service to a revision

As you noticed that the Knative services always routes the traffic to the **latest** revision of the service deployment, thats because of the **runLatest** attribute in [service resource file](#).

Let us now make the greeter service use earlier revision **greeter-00001**.



You can use the get [\[show-knative-revisions\]](#) command to find the available revisions for the greeter service.

## Service pinned to revision 1

```
apiVersion: serving.knative.dev/v1alpha1
kind: Service
metadata:
  name: greeter
spec:
  pinned: ①
  revisionName: greeter-00001
  configuration:
    revisionTemplate:
      spec:
        container:
          image: dev.local/rhdevelopers/greeter:0.0.1
```

### service-pinned-rev1.yaml

- ① The **pinned** attribute in service resource file will make Knative use the revision specified in the **revisionName** attribute.

Let redeploy the greeter service to be pinned to revision *greeter-00001*:

```
cd $TUTORIAL_HOME/knative
```

```
oc apply -n knativetutorial -f service-pinned-rev1.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f service-pinned-rev1.yaml
```

[Invoking Service](#) will now show an output like **Hi greeter** ⇒ **greeter-00001-deployment-5d696cc6c8-m65s5: 3**.

## Cleanup

```
oc -n knativetutorial delete services.serving.knative.dev greeter
```

**(OR)**

```
kubectl -n knativetutorial delete services.serving.knative.dev greeter
```

# Configurations and Routes

At the end of this chapter you will be able to understand and know how to :

- Deploying Knative configurations and routes separately
- Distributing traffic between revisions of the service

## Prerequisite

The following checks ensures that each chapter exercises are done with right environment settings.

- Make sure to be connected to minishift docker daemon and using right openshift client.

```
eval $(minishift docker-env) && eval $(minishift oc-env)
```

- Make sure right OpenShift and kubernetes client are used.

```
# oc v3.11.0+0cbc58b kubernetes v1.11.0+d4cacc0
oc version
```

(OR)

```
# kubernetes v1.11+
kubectl version
```

- Make sure right maven version is used (only for Java exercises)

```
# right maven version Apache Maven 3.5.x or above
./mvnw --version
```

- Make sure to be on **knativetutorial** OpenShift project/kubernetes namespace

```
# verify to return knativetutorial
oc project -q
```

If you are not on **knativetutorial** project, use the command **oc project knativetutorial** to change to **knativetutorial** project

## Build Containers

- Navigate to work folder

```
cd $TUTORIAL_HOME/work
```

- Clone the application source code

```
git clone https://github.com/redhat-developer-demos/knative-tutorial-greeter.git
```

- Build the application and container image

```
cd knative-tutorial-greeter/java/springboot
./mvnw clean package
docker build -t dev.local/rhdevelopers/greeter:0.0.1 .
```

- Check if images are built and available:

```
docker images | grep dev.local
```

The above command should return something like as shown below:

```
# dev.local/rhdevelopers/greeter 0.0.1 6aa8771da8dd 2 hours ago
456MB
```



You can also use the script `buildImage.sh` that will build source and create the container image

## Deploy Configuration

In previous chapter we saw how we can deploy service using holistic service resource file, in this chapter we will see how to deploy the service using configurations and resource files.

Navigate to folder `$TUTORIAL_HOME/02-configs-and-routes`:

### Deploy Configuration revision 1

The following snippet how a Knative configuration resource YAML will look like:

```

apiVersion: serving.knative.dev/v1alpha1
kind: Configuration
metadata:
  name: greeter
spec:
  revisionTemplate:
    metadata:
      labels:
        app: greeter
    spec:
      container:
        image: dev.local/rhdevelopers/greeter:0.0.1 ①

```

### configuration-rev1.yaml

① It is very important that the image is a fully qualified name docker image name with tag. For more details on this [Question 2 of FAQ](#)

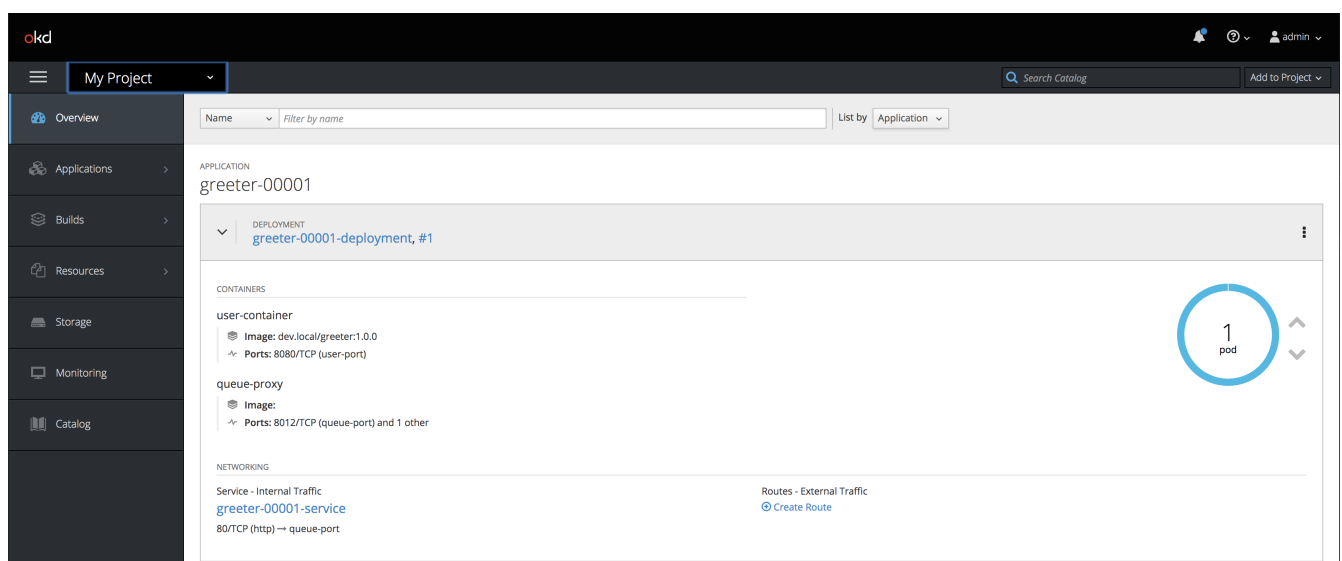
The service could be deployed using the command:

```
oc apply -n knativetutorial -f config/configuration-rev1.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f config/configuration-rev1.yaml
```

After successful deployment of the service we should see a kubernetes deployment called **greeter-00001-deployment** available.



## Invoke Service

```
INGRESSGATEWAY=istio-ingressgateway
IP_ADDRESS="$(minishift ip):$(kubectl get svc $INGRESSGATEWAY --namespace istio-system
--output 'jsonpath={.spec.ports[?(@.port==80)].nodePort}')"

```

```
curl -H "Host: greeter.knativetutorial.example.com" $IP_ADDRESS

```

(OR)

```
http $IP_ADDRESS 'Host: greeter.knativetutorial.example.com'

```

The command will return HTTP 404 as there are no routes deployed yet. Let us now deploy a route

## Deploy Route

Let us now deploy a default route that will route all the traffic to service deployed via configuration `greeter`

```
apiVersion: serving.knative.dev/v1alpha1
kind: Route
metadata:
  name: greeter
spec:
  traffic:
    - configurationName: greeter
      percent: 100

```

**route\_default.yaml**

```
oc apply -n knativetutorial -f route/route_default.yaml

```

(OR)

```
kubectl apply -n knativetutorial -f route/route_default.yaml

```

[Invoking Service](#) now should return a response like **Hi greeter** ⇒ **greeter-00001-deployment-5d696cc6c8-m65s5: 1**



Sometimes the response might not be returned immediately especially when the pod is coming up from dormant state, at those times try giving request again



# See what you have deployed

The holistic service base deployment strategy that we did now will create the following resource objects:

## service

```
# get a Knative Service (short name ksvc) called greeter
oc -n knativetutorial get services.serving.knative.dev greeter
```

(OR)

```
kubectl -n knativetutorial get services.serving.knative.dev greeter
```

The command will throw an error as there are no knative services deployed as part of configuration+routes based deployment strategy.

## configuration

```
# get a Knative configuration called greeter
oc -n knativetutorial get configurations.serving.knative.dev greeter
```

(OR)

```
kubectl -n knativetutorial get configurations.serving.knative.dev greeter
```

## routes

```
# get a Knative routes called greeter
oc -n knativetutorial get routes.serving.knative.dev greeter
```

(OR)

```
kubectl -n knativetutorial get routes.serving.knative.dev greeter
```

When the service was invoked with `curl -H "Host: greeter.knativetutorial.example.com" $IP_ADDRESS`, you noticed that we added a **Host** header to the curl with value `greeter.knativetutorial.example.com`, this FQDN is automatically assigned to your Knative service by the Knative Routes. It follows a format `<service-name>.<namespace>.<domain-suffix>`.



- The domain suffix in this case *example.com* is configurable via the config map **config-domain** of **knative-serving** namespace.

## revisions

```
# get a Knative revisions called greeter you will see only one like greeter-00001
oc -n knativetutorial get revisions.serving.knative.dev
```

(OR)

```
kubectl -n knativetutorial get revisions.serving.knative.dev
```



- add `-oyaml` to the commands above to see more details

## Deploy new Revision Service

For the new revision of the service we will add an environment variable to the service deployment. As Knative follows 12-Factor application principles, any new configuration change will trigger new deployment, in this case it is the new environment variable.

### Deploy configuration revision 2

```
apiVersion: serving.knative.dev/v1alpha1
kind: Configuration
metadata:
  name: greeter
spec:
  revisionTemplate:
    metadata:
      labels:
        app: greeter
    spec:
      container:
        image: dev.local/rhdevelopers/greeter:0.0.1
        env: ①
        - name: MESSAGE_PREFIX
          value: Hello
```

**configuration-rev2.yaml**

- ① Adding an environment variable that will be used a message prefix

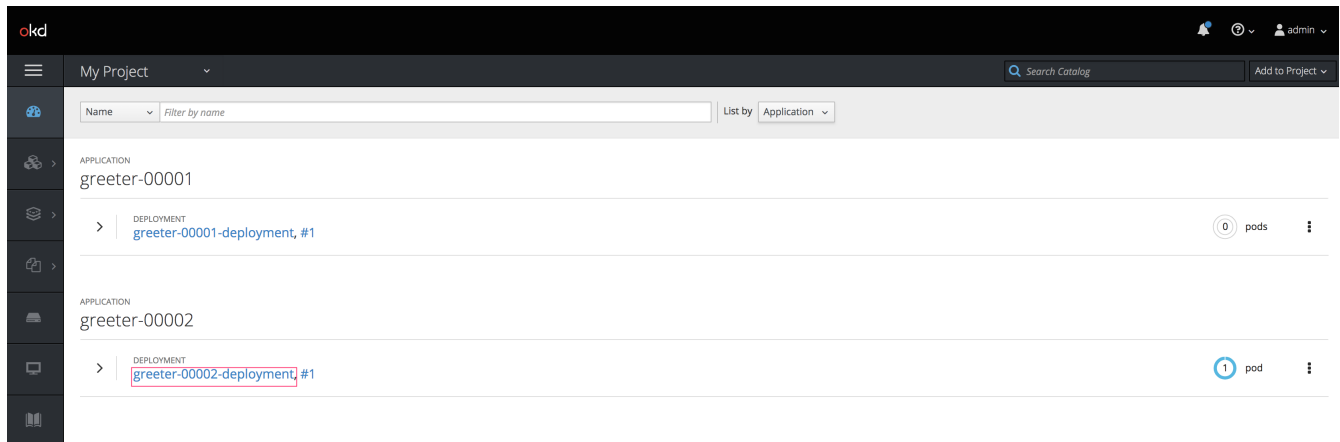
Let us deploy the new revision using the command:

```
oc apply -n knativetutorial -f config/configuration-rev2.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f config/configuration-rev2.yaml
```

After successful deployment of the service we should see a kubernetes deployment called **greeter-00002-deployment** available in the OpenShift dashboard:



Now running the [command](#) will show two revisions namely **greeter-00001** and **greeter-0002**.

[Invoking Service](#) will now show an output like **Hello greeter** ⇒ **greeter-00002-deployment-8d9984dc8-rgzx6: 2**, where *Hello* is the value that we configured via environment variable in the Knative service resource file.

## Distributing traffic

When deploying services with service object based approach, the service deployment will take care of traffic distribution either 100% to latest or pinned revision. In this section we will explore how we can deploy routes to distribute traffic between deployed revisions.



For the sake of clarity we will call greeter-00001 as **revision 1** and greeter-00002 as **revision 2**

Before we start to apply routes, let us open a new terminal and run the following command `$TUTORIAL_HOME/02-configs-and-routes/bin/call.sh`, this command will keep sending requests to greeter service in every two seconds which allows us to monitor the changes to responses as we keep applying the route that will distribute the traffic between available two revisions.

### All traffic to revision 1

```
oc apply -n knativetutorial -f route/route_all_rev1.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f route/route_all_rev1.yaml
```

You will notice the output on your other terminal to be something like **Hi greeter** ⇒ **greeter-00001-**

**deployment-5d696cc6c8-m65s5: 34.**

## All traffic to revision 2

```
oc apply -n knativetutorial -f route/route_all_rev2.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f route/route_all_rev2.yaml
```

You will notice the output on your other terminal to be something like **Hello greeter** ⇒ **greeter-00002-deployment-8d9984dc8-rgzx6: 13.**

## 50-50 between revision 1 and revision 2

```
oc apply -n knativetutorial -f route/route_rev1-50_rev2-50.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f route/route_rev1-50_rev2-50.yaml
```

You will notice the output will be mix of responses like **Hi greeter** ⇒ **greeter-00001-deployment-5d696cc6c8-m65s5: 11** and **Hello greeter** ⇒ **greeter-00002-deployment-8d9984dc8-rgzx6: 10** approximately distributed 50% between each.

## 75-25 between revision 1 and revision 2

```
oc apply -n knativetutorial -f route/route_rev1-75_rev2-25.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f route/route_rev1-75_rev2-25.yaml
```

You will notice the output will be mix of responses like **Hi greeter** ⇒ **greeter-00001-deployment-5d696cc6c8-m65s5: 6** and **Hello greeter** ⇒ **greeter-00002-deployment-8d9984dc8-rgzx6: 7**, with more requests responded by revision 1 approximately distributed 75% to 25 % between revision 1 and revision 2.

## 10-90 between revision 1 and revision 2

```
oc apply -n knativetutorial -f route/route_rev1-10_rev2-90.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f route/route_rev1-10_rev2-90.yaml
```

You will notice there will be a mix of responses like **Hi greeter** ⇒ **greeter-00001-deployment-5d696cc6c8-m65s5: 4** and **Hello greeter** ⇒ **greeter-00002-deployment-8d9984dc8-rgzx6: 5**, with more requests responded by revision 2 approximately 10% to 90% between revision 1 and revision 2.



In response texts e.g **Hi greeter** ⇒ **greeter-00001-deployment-5d696cc6c8-m65s5: 4** the numbers at the end of the responses shows a count of how many requests have been handled by the service in this example it is 4. Also note that the output count number may vary according to number of requests that you might have issued.

## Cleanup

```
oc -n knativetutorial delete configurations.serving.knative.dev greeter
oc -n knativetutorial delete routes.serving.knative.dev greeter
```

(OR)

```
kubectl -n knativetutorial delete configurations.serving.knative.dev greeter
kubectl -n knativetutorial delete routes.serving.knative.dev greeter
```

# Scaling

At the end of this chapter you will be able to:

- Configuring the scale to zero time period
- Configuring auto scaling
- Understanding types of autoscaling strategies
- Enabling concurrency based auto scaling based concurrency
- Configuring minimum number of replicas for a service i.e. no scale to zero below the replica count

## Prerequisite

The following checks ensures that each chapter exercises are done with right environment settings.

- Make sure to be connected to minishift docker daemon and using right openshift client.

```
eval $(minishift docker-env) && eval $(minishift oc-env)
```

- Make sure right OpenShift and kubernetes client are used.

```
# oc v3.11.0+0cbc58b kubernetes v1.11.0+d4cacc0
oc version
```

(OR)

```
# kubernetes v1.11+
kubectl version
```

- Make sure right maven version is used (only for Java exercises)

```
# right maven version Apache Maven 3.5.x or above
./mvnw --version
```

- Make sure to be on **knativetutorial** OpenShift project/kubernetes namespace

```
# verify to return knativetutorial
oc project -q
```

If you are not on **knativetutorial** project, use the command **oc project knativetutorial** to change to **knativetutorial** project

# Build Containers

- Navigate to work folder

```
cd $TUTORIAL_HOME/work
```

- Clone the application source code

```
git clone https://github.com/redhat-developer-demos/knative-tutorial-greeter.git
```

- Build the application and container image

```
cd knative-tutorial-greeter/java/springboot
./mvnw clean package
docker build -t dev.local/rhdevelopers/greeter:0.0.1 .
```

- Check if images are built and available:

```
docker images | grep dev.local
```

The above command should return something like as shown below:

```
# dev.local/rhdevelopers/greeter 0.0.1 6aa8771da8dd 2 hours ago
456MB
```



You can also use the script `buildImage.sh` that will build source and create the container image

## Deploy Service

The following snippet shows how a Knative service YAML will look like:

```

apiVersion: serving.knative.dev/v1alpha1
kind: Service
metadata:
  name: greeter
spec:
  runLatest: ①
  configuration:
    revisionTemplate:
      spec:
        container:
          image: dev.local/rhdevelopers/greeter:0.0.1 ②

```

### service.yaml

- ① Makes Knative to always run the latest revision of the deployment
- ② It is very important that the image is a fully qualified name docker image name with tag. For more details on this [Question 2 of FAQ](#)

```
cd $TUTORIAL_HOME/01-basics/knative
```

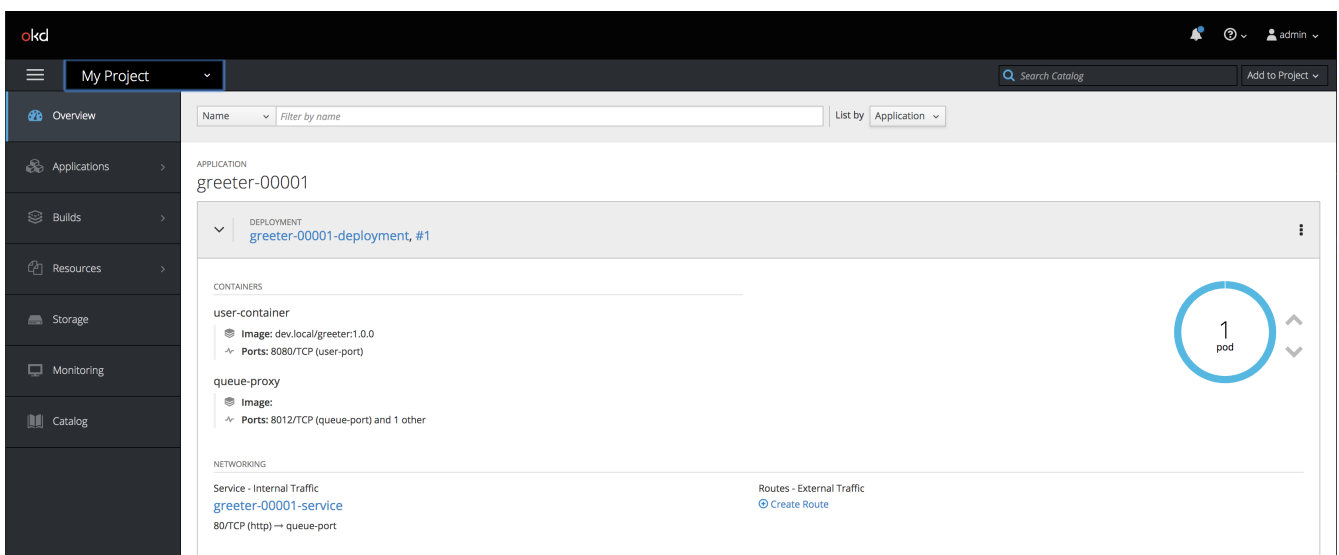
The service could be deployed using the command:

```
oc apply -n knativetutorial -f service.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f service.yaml
```

After successful deployment of the service we should see a kubernetes deployment called **greeter-00001-deployment** available.





# Invoke Service

```
INGRESSGATEWAY=istio-ingressgateway
IP_ADDRESS="$(minishift ip):$(kubectl get svc $INGRESSGATEWAY --namespace istio-system --output 'jsonpath={.spec.ports[?(@.port==80)].nodePort}')"

```

```
curl -H "Host: greeter.knativetutorial.example.com" $IP_ADDRESS

```

(OR)

```
http $IP_ADDRESS 'Host: greeter.knativetutorial.example.com'

```

The last curl command should return a response like **Hi greeter** ⇒ **greeter-00001-deployment-5d696cc6c8-m65s5: 1**

Check [deployed Knative resources](#) for more details of what Knative objects and resources have been created with the above service deployment.

## Configure scale to zero

Assuming that [Greeter service](#) has been deployed, if the service has been terminated, then let us [invoke the service](#) to make service available for request.

### Calculating scale to zero time period

Knative autoscaler uses the config map **config-autoscaler** of **knative-serving** namespace for autoscaling related properties. The value of the attribute **stable-window** and **scale-to-zero-grace-period** decides how long Knative autoscaler will wait before terminating the inactive revision pod.

```
stableWindow='oc -n knative-serving get configmaps config-autoscaler -o yaml |
yq r - data.stable-window'
scaleToZeroGracePeriod='oc -n knative-serving get configmaps config-autoscaler -o yaml |
yq r - data.scale-to-zero-grace-period'

```

(OR)

```
stableWindow='kubectl -n knative-serving get configmaps config-autoscaler -o yaml |
yq r - data.stable-window'
scaleToZeroGracePeriod='kubectl -n knative-serving get configmaps config-autoscaler -o yaml |
yq r - data.scale-to-zero-grace-period'

```

**termination period time in seconds = stableWindow + scaleToZeroGracePeriod**



- `scale-to-zero-grace-period` is dynamic parameter i.e. the value is immediately effected after updating the config map
- the minimum value that can be set for `scale-to-zero-grace-period` is 30 seconds

## Observing default scale down

For easier observation let us open a new terminal and run the following command,

```
watch 'oc get pods -n knativetutorial'
```



The scaling up and down can also be watched on OpenShift dashboard by navigating to knativetutorial project

Checking default values

```
# should return 60s
echo $stableWindow
# should return 30s
echo $scaleToZeroGracePeriod
```

By default the **scale-to-zero-grace-period** is **30s**, and the **stable-window** is **60s**, firing the request to the greeter service will bring up the pod (if its already terminated) to serve the request. Leaving it without any further requests, it will automatically scale to zero **1 min 30 secs**([Compute scale to zero grace period](#)).

## Update scale down to 1 minute

Let us now update **scale-to-zero-grace-period** to **1m** and leave the **stable-window** to default **60s**.

```
cd $TUTORIAL_HOME/03-scaling/knative
```

```
oc -n knative-serving get cm config-autoscaler -o yaml | yq w - -s configure-scaling-to-1m.yaml | oc apply -f -
```

(OR)

```
kubectl -n knative-serving get cm config-autoscaler -o yaml | yq w - -s configure-scaling-to-1m.yaml | kubectl apply -f -
```

Verifying the `scale-to-zero-grace-period` value, which should return **1m**:

```
oc -n knative-serving get configmap config-autoscaler -o yaml \
| yq r - data.scale-to-zero-grace-period
```

(OR)

```
kubectl -n knative-serving get configmap config-autoscaler -o yaml \
| yq r - data.scale-to-zero-grace-period
```

Now [firing the request](#) to the greeter service will bring up the pod to serve the request. Leaving it without any further requests, it will automatically scale to zero **2 mins**([Compute scale to zero grace period](#))

## Update scale down to 2 minute

Let us now update **scale-to-zero-grace-period** to **2m** and leave the **stable-window** to default **60s**.

```
oc -n knative-serving get cm config-autoscaler -o yaml | yq w - -s configure-scaling-to-2m.yaml | oc apply -f -
```

(OR)

```
kubectl -n knative-serving get cm config-autoscaler -o yaml | yq w - -s configure-scaling-to-2m.yaml | kubectl apply -f -
```

Verifying the **scale-to-zero-grace-period** value, which should return **2m**:

```
# should return 2m
oc -n knative-serving get configmap config-autoscaler -o yaml \
| yq r - data.scale-to-zero-grace-period
```

(OR)

```
kubectl -n knative-serving get configmap config-autoscaler -o yaml \
| yq r - data.scale-to-zero-grace-period
```

Now [firing the request](#) to the greeter service will bring up the pod to serve the request and if we leave the service without any further requests, it will automatically scale to zero in **3 mins**([Compute scale to zero grace period](#)).

## Reset to defaults

Let us revert the **scale-to-zero-grace-period** to defaults:

```
oc -n knative-serving get cm config-autoscaler -o yaml | yq w - -s configure-scaling-to-30s.yaml | oc apply -f -
```

(OR)

```
kubectl -n knative-serving get cm config-autoscaler -o yaml | yq w - -s configure-scaling-to-30s.yaml | kubectl apply -f -
```

Verifying the **scale-to-zero-grace-period** value, which should return **30s**:

```
# should return 2m
oc -n knative-serving get configmap config-autoscaler -o yaml \
| yq r - data.scale-to-zero-grace-period
```

(OR)

```
kubectl -n knative-serving get configmap config-autoscaler -o yaml \
| yq r - data.scale-to-zero-grace-period
```

For better clarity and understanding let us **clean up** the deployed knative resources before going to next section.

## Auto Scaling

By default Knative Serving allows 100 concurrency pods, you can view the setting **container-concurrency-target-default** in the configmap **config-autoscaler** of **knative-serving** namespace.

For this exercise let us make our service handle only **10** concurrent requests, this will make the Knative-Serving to scale the pods to handle the requests more than 10 requests

### Service with concurrency of 10 requests

```

apiVersion: serving.knative.dev/v1alpha1
kind: Service
metadata:
  name: greeter
spec:
  runLatest:
    configuration:
      revisionTemplate:
        metadata:
          annotations:
            autoscaling.knative.dev/target: "10" ①
        spec:
          container:
            image: dev.local/rhdevelopers/greeter:0.0.1
            livenessProbe:
              httpGet:
                path: /healthz
            readinessProbe:
              httpGet:
                path: /healthz

```

#### service-10.yaml

- ① Will allow each service pod to handle max of 10 in-flight requests per pod before automatically scaling to new pod(s)

## Deploy service

```
oc apply -n knativetutorial -f service-10.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f service-10.yaml
```

## Invoke Service

```

INGRESSGATEWAY=istio-ingressgateway
IP_ADDRESS="$(minishift ip):$(kubectl get svc $INGRESSGATEWAY --namespace istio-system
--output 'jsonpath={.spec.ports[?(@.port==80)].nodePort}')"

```

Open a new terminal and run the following command:

```
watch oc get pods -n knativetutorial
```

(OR)

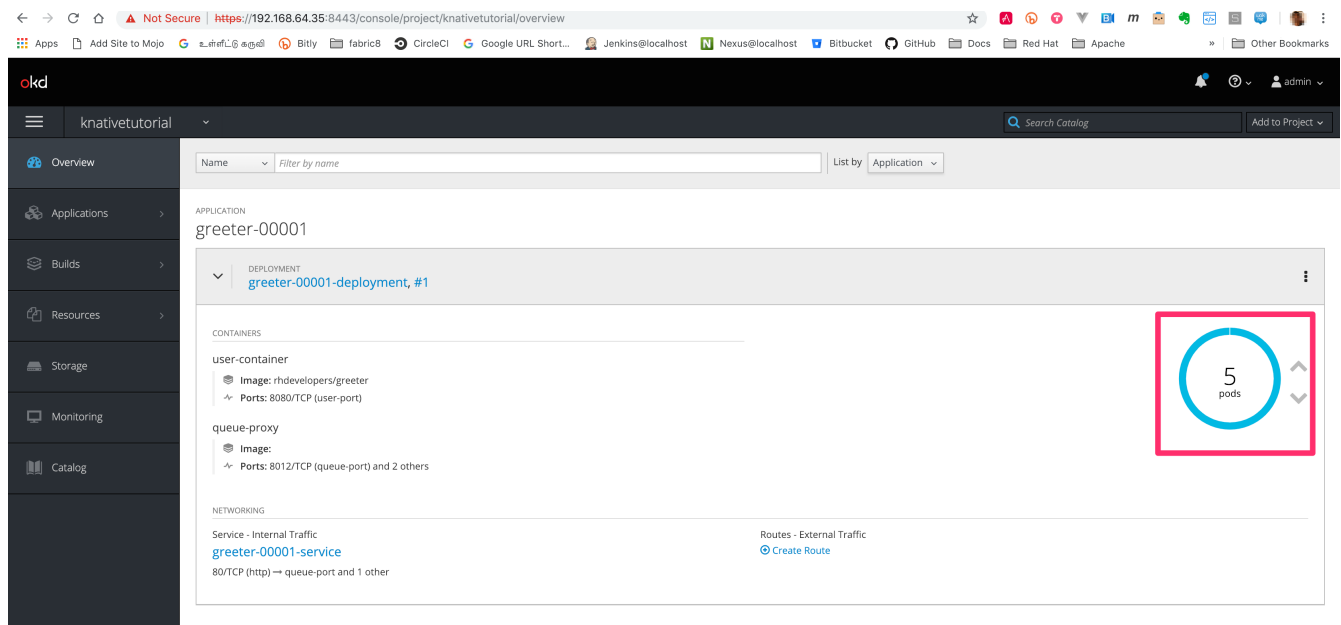
```
watch kubectl get pods -n knativetutorial
```

## Load the service

We will now send some load to the greeter service. The command below sends 50 concurrent requests (-c 50) for next 10s (-z 10s) with no connection timeout(-t 0)

```
hey -z 10s -c 50 -t 0 \  
-host "greeter.knativetutorial.example.com" \  
"http://${IP_ADDRESS}"
```

After successful run of the load test, you will notice the number of greeter service pods would have automatically scaled to 5. You can also view the same via OpenShift dashboard as shown below:



The autoscale pods is computed using the formula:

$$\text{totalPodsToScale} = \text{total number of inflight requests} / \text{average concurrency target}$$

With this current setting of **average concurrency target=10**(`autoscaling.knative.dev/target: "10"`) and **total number of inflight requests=50** , you will see Knative automatically scales the greeter services to **50/10 = 5 pods**.

For more clarity and understanding let us [clean up](#) existing deployments before proceeding to next section.

## Minimum Scale

In real world scenarios your service might need to handle sudden spikes in requests. Knative by

default set the replica count of the deployed service(s) to be one, but this default setting can be configured via an service annotations **autoscaling.knative.dev/minScale**.

The following example shows how to make knative services by default create services with replica count of two.

## Deploy service

```
apiVersion: serving.knative.dev/v1alpha1
kind: Service
metadata:
  name: greeter
spec:
  runLatest:
    configuration:
      revisionTemplate:
        metadata:
          annotations:
            autoscaling.knative.dev/target: "10" ①
            autoscaling.knative.dev/minScale: "2" ②
        spec:
          container:
            image: dev.local/rhdevelopers/greeter:0.0.1
            livenessProbe:
              httpGet:
                path: /healthz
            readinessProbe:
              httpGet:
                path: /healthz
```

### service-min-scale.yaml

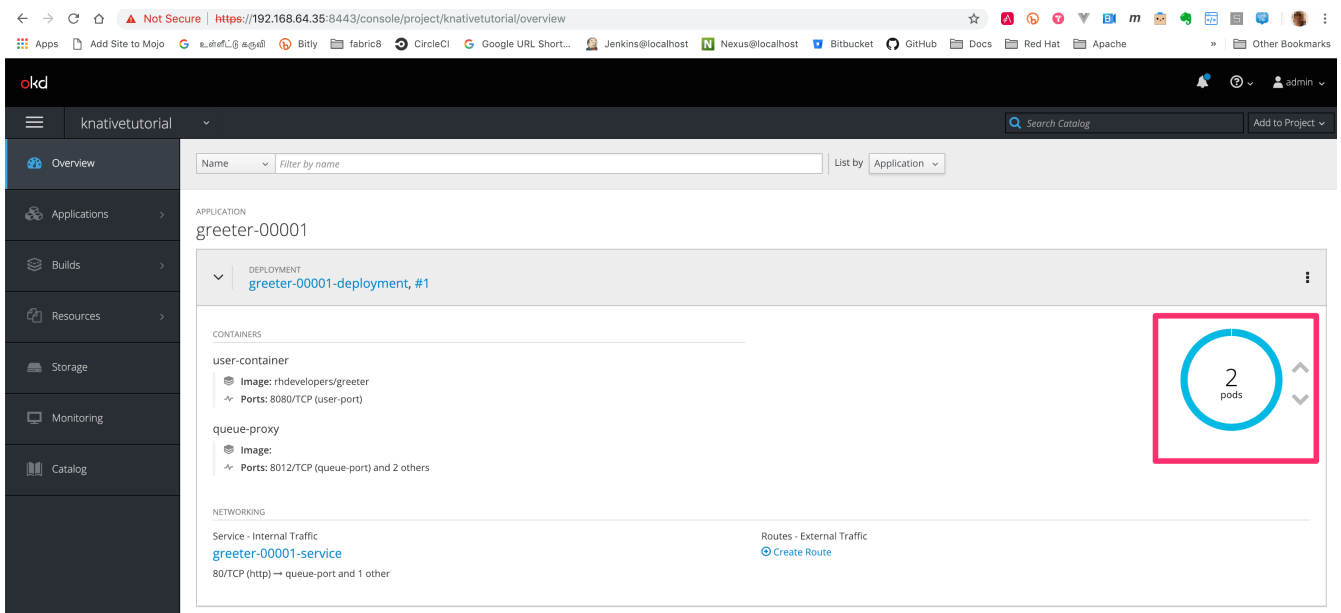
- ① Will allow each service pod to handle max of 10 in-flight requests per pod before automatically scaling to new pod(s)
- ② Every new deployment of this service will minimum have two pods

```
oc apply -n knativetutorial -f service-min-scale.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f service-min-scale.yaml
```

After successful deployment of the service we should see a kubernetes deployment called **greeter-00001-deployment** with **two** pods readily available.



## OpenShift Dashboard Knative Tutorial project

Open a new terminal and run the following command :

```
watch oc get pods -n knativetutorial
```

(OR)

```
watch kubectl get pods -n knativetutorial
```

Let us [send some load to the service](#) to trigger autoscaling.

When all requests are done and if we are beyond the [scale-to-zero-grace-period](#), we will notice that Knative has terminated only 3 out 5 pods, this is because we have configured Knative to always run two pods via the annotation [autoscaling.knative.dev/minScale: "2"](#)

## Cleanup

```
oc -n knativetutorial delete services.serving.knative.dev greeter
```

(OR)

```
kubectl -n knativetutorial delete services.serving.knative.dev greeter
```



# Build

At the end of this chapter you will be able to:

- Define build
- Define build templates
- How to make service use builds for building container images a.k.a source-to-url

## Prerequisite

The following checks ensures that each chapter exercises are done with right environment settings.

- Make sure to be connected to minishift docker daemon and using right openshift client.

```
eval $(minishift docker-env) && eval $(minishift oc-env)
```

- Make sure right OpenShift and kubernetes client are used.

```
# oc v3.11.0+0cbc58b kubernetes v1.11.0+d4cacc0
oc version
```

(OR)

```
# kubernetes v1.11+
kubectl version
```

- Make sure right maven version is used (only for Java exercises)

```
# right maven version Apache Maven 3.5.x or above
./mvnw --version
```

- Make sure to be on **knativetutorial** OpenShift project/kubernetes namespace

```
# verify to return knativetutorial
oc project -q
```

If you are not on **knativetutorial** project, use the command **oc project knativetutorial** to change to **knativetutorial** project



It's required that you have a container registry account which is needed during the exercises to push the built container images.

You can get one for free at [Docker Hub](#)

As the Knative build and build template spec needs some customizations, such as creating a secret with your container registry credentials; we need to run the following commands to create the needed Kubernetes and Knative resources with all your customizations. All the generated files will be placed in `$TUTORIAL_HOME/04-build/knative`.

## Generate Docker Secret

A Docker Kubernetes secret is required for the build to authenticate with the Docker registry so that it can push the built container image to the registry.

```
jsonnet --ext-str user='<your-docker-repo-userid>' \  
  --ext-str password='<your-docker-repo-user-password>' \  
  docker-secret.jsonnet \  
  | yq r - | tee ../docker-secret.yaml
```

(e.g.)

```
jsonnet --ext-str user='demo' \  
  --ext-str password='demopassword' \  
  docker-secret.jsonnet \  
  | yq r - | tee ../docker-secret.yaml
```

The above command will generate a Kubernetes secret as shown below:

```
apiVersion: v1  
kind: Secret  
metadata:  
  annotations:  
    build.knative.dev/docker-0: https://index.docker.io/v1/ ①  
  name: basic-user-pass  
stringData:  
  password: demopassword ②  
  username: demo ③  
type: kubernetes.io/basic-auth
```

### docker-secret.yaml

- ① The username and password that are specified as part of this secret will be used for Docker container registry available at <https://index.docker.io/v1/>
- ② The username to use when authenticating with <https://index.docker.io/v1/>

③ The password to use when authenticating with <https://index.docker.io/v1/>

For more details on why we need to do this please check [here](#).

## Generate Knative build spec

The knative build spec is used to create and run the knative builds. You can find more details on Knative build [here](#).

```
jsonnet --ext-str contextDir='<the build context directory>' \  
  --ext-str image='<your-docker-image-name>' \  
  docker-build.jsonnet \  
  | yq r - | tee ../docker-build.yaml
```

(e.g.)

```
jsonnet --ext-str contextDir='java/springboot' \  
  --ext-str image='docker.io/demo/event-greeter:0.0.1' \  
  docker-build.jsonnet \  
  | yq r - | tee ../docker-build.yaml
```

The above command will generate a knative build specification as shown below:

```

apiVersion: build.knative.dev/v1alpha1
kind: Build
metadata:
  name: docker-build
spec:
  serviceAccountName: build-bot
  source:
    git:
      revision: master
      url: https://github.com/redhat-developer-demos/knative-tutorial-event-
greeter.git ①
  steps: ②
  - args:
    - --context=/workspace/java/springboot ③
    - --dockerfile=/workspace/java/springboot/Dockerfile ④
    - --destination=docker.io/demo/event-greeter:0.0.1 ⑤
    env:
    - name: DOCKER_CONFIG
      value: /builder/home/.docker
    image: gcr.io/kaniko-project/executor
    name: docker-push
    volumeMounts:
    - mountPath: /builder/home/.m2
      name: m2-cache
    - mountPath: /cache
      name: kaniko-cache
    workingDir: /workspace/java/springboot
  timeout: 20m
  volumes:
  - name: m2-cache
    persistentVolumeClaim:
      claimName: m2-cache
  - name: kaniko-cache
    persistentVolumeClaim:
      claimName: kaniko-cache

```

### docker-build.yaml

- ① The Git project repository which will be built as part of this build
- ② Each build can have one or more steps, in this case its just one step
- ③ Set from `containerDir` parameter—the sub directory relative to the root of the project from where the build will be run. In this example it is the `java/springboot` —
- ④ The *Dockerfile* is assumed to be in the root of the context directory
- ⑤ Set from `image` parameter— the fully qualified container image that will be pushed to the registry—

# Generate Knative service

Run the following command to create the knative service that will use the image created by the `docker-build`.

```
jsonnet --ext-str image='<your-docker-image-name>' \  
  service.jsonnet \  
  | yq r - | tee ../service-build.yaml
```

```
jsonnet --ext-str image='docker.io/demo/event-greeter:0.0.1' \  
  service.jsonnet \  
  | yq r - | tee ../service-build.yaml
```

```
apiVersion: serving.knative.dev/v1alpha1  
kind: Service  
metadata:  
  name: event-greeter  
spec:  
  runLatest:  
    configuration:  
      revisionTemplate:  
        metadata:  
          labels:  
            app: event-greeter  
        spec:  
          container:  
            image: docker.io/demo/event-greeter:0.0.1 ①
```

## service-build.yaml

① Set from `image` parameter — the fully qualified container image that will be pushed to the registry as part of the build —

## Apply pre-requisite resources

Before we create the build and service, we need to create a few kubernetes resources that will be used by the build.

Run the following commands to create the pre-requisite resources:

```
oc apply -n knativetutorial -f docker-secret.yaml ①  
oc apply -n knativetutorial -f build-sa.yaml ②  
oc apply -n knativetutorial -f m2-pvc.yaml ③  
oc apply -n knativetutorial -f kaniko-pvc.yaml ④
```

(OR)

```
kubectl apply -n knativetutorial -f docker-secret.yaml
kubectl apply -n knativetutorial -f build-sa.yaml
kubectl apply -n knativetutorial -f m2-pvc.yaml
kubectl apply -n knativetutorial -f kaniko-pvc.yaml
```

- ① The kubernetes secret called **basic-auth-pass**, check [here](#) for more details
- ② The kubernetes service account called **build-bot** that will run the build. Only **build-bot** has access to the secret **basic-auth-pass**. Check [here](#) for more details
- ③ the kubernetes persistence volume claim used to cache maven artifacts, to make subsequent builds faster
- ④ the kubernetes persistence volume claim that will be used to cache docker base images to make subsequent [Kaniko](#) docker builds faster

## Create build

```
oc apply -n knativetutorial -f docker-build.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f docker-build.yaml
```

## See what you have deployed

### build

```
# get a knative builds called event-greeter
oc -n knativetutorial get builds.build.knative.dev event-greeter
```

(OR)

```
kubectl -n knativetutorial get builds.build.knative.dev event-greeter
```

## Watching logs

The Knative build creates an init container for each build step. The init containers will be named like **build-step-<step-name>** e.g. if the build step is named "foo" then there will be an init container named "build-step-foo" added to the build pod.

The following command can be used to watch the log of a specific container within a pod,

```
oc logs -n knativetutorial <pod-name> -c <build-init-container-name>
```

(OR)

```
kubectl logs -n knativetutorial <pod-name> -c <build-init-container-name>
```

(e.g.)

To watch the container logs of container `docker-push` within the build pod called `docker-build-pod-2a271e`:

```
oc logs -n knativetutorial docker-build-pod-2a271e -c build-step-docker-push
```

(OR)

```
kubectl logs -n knativetutorial docker-build-pod-2a271e -c build-step-docker-push
```



- You can use the following command to list names of all the init containers within a pod .

(e.g.) the following command lists all the init container names of the pod `docker-build-pod-2986f2`

```
oc -n knativetutorial get pods docker-build-pod-2986f2 -o yaml | yq r - spec.initContainers[*].name
```

(OR)

```
kubectl -n knativetutorial get pods docker-build-pod-2986f2 -o yaml | yq r - spec.initContainers[*].name
```

- Using `stern` you can watch logs of all the container of the build pod using the command `stern docker-build-pod`

## Checking build status

A successful build will have the container image pushed to the container registry, going by the examples above you should the image `event-greeter:0.0.1`` in your container registry.

You can check the status of build as shown below:

```
$ oc -n knativetutorial get pods
```

(OR)

```
$ kubectl -n knativetutorial get pods
```

The above command should return an output as shown below:

NAME	READY	STATUS	RESTARTS	AGE
docker-build-pod-ecc235	0/1	Completed	0	10m

With this you should also be able to do `docker pull docker.io/demo/event-greeter:0.0.1` successfully.

## Deploy service using Build

With a successful build you are ready to deploy the service, run the following commands to deploy the service:

```
oc apply -n knativetutorial -f service-build.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f service-build.yaml
```

After successful deployment of the service we should see a kubernetes deployment called `event-greeter-00001-deployment` available in the OpenShift dashboard:

The screenshot shows the OpenShift Application Console interface. The top navigation bar includes the 'okd' logo and the 'Application Console' dropdown. The main content area is titled 'knativetutorial' and shows a list of applications. The 'event-greeter' application is selected, and its details are displayed. The 'DEPLOYMENT' section shows 'event-greeter-00001-deployment, #1'. The 'CONTAINERS' section lists two containers: 'user-container' and 'queue-proxy'. The 'user-container' is highlighted, showing its image as 'kameshsampath/event-greeter' and its ports as '8080/TCP (user-port)'. The 'queue-proxy' container is also listed. The 'NETWORKING' section shows the service 'event-greeter-00001-service' and its ports. A red box highlights a circular icon with the number '1' and the text 'pod', indicating the number of pods running for the deployment.



## Invoke service

```
INGRESSGATEWAY=istio-ingressgateway
IP_ADDRESS="$(minishift ip):$(kubectl get svc $INGRESSGATEWAY --namespace istio-system
--output 'jsonpath={.spec.ports[?(@.port==80)].nodePort}')"

```

```
curl -H "Host: event-greeter.knativetutorial.example.com" $IP_ADDRESS \
-H "Content-Type: application/json" \
-d '{"message": "test message"}'
```

(OR)

```
http POST $IP_ADDRESS \
'Host: event-greeter.knativetutorial.example.com' \
message="test message"
```

The above command should return an response like:

```
{
  "host": "Event greeter => 'unknown' : 1 ",
  "message": "test message"
}
```

## Cleanup

Before proceeding next chapter let us delete the build and service objects that were created.

```
oc -n knativetutorial delete services.serving.knative.dev event-greeter
oc -n knativetutorial delete builds.build.knative.dev docker-build
```

(OR)

```
kubectl -n knativetutorial delete services.serving.knative.dev event-greeter
kubectl -n knativetutorial delete builds.build.knative.dev docker-build
```

# Build Template

The raw builds like what we did in [previous section](#) has some disadvantages

- No reusability
- Hard coding of build parameters
- Inability to make builds and services in synchronous i.e. services are deployed only after the related build is complete

One of the solution to these problems is by using **Build Templates**.

Like the build spec the knative build template spec and knative service spec also needs some customizations before apply them.

## Prerequisite

The following checks ensures that each chapter exercises are done with right environment settings.

- Make sure to be connected to minishift docker daemon and using right openshift client.

```
eval $(minishift docker-env) && eval $(minishift oc-env)
```

- Make sure right OpenShift and kubernetes client are used.

```
# oc v3.11.0+0cbc58b kubernetes v1.11.0+d4cacc0  
oc version
```

(OR)

```
# kubernetes v1.11+  
kubectl version
```

- Make sure right maven version is used (only for Java exercises)

```
# right maven version Apache Maven 3.5.x or above  
./mvnw --version
```

- Make sure to be on **knativetutorial** OpenShift project/kubernetes namespace

```
# verify to return knativetutorial  
oc project -q
```

If you are not on **knativetutorial** project, use the command **oc project knativetutorial** to change

to [knativetutorial](#) project



- Please ensure that you have [registry account](#) created.
- Pre-requisite [kubernetes resources](#) have been created and applied.

## Generate build template

```
[source,bash,linenums]
----
jsonnet maven-build-template.jsonnet \
  | yq r - | tee ../maven-build-template.yaml
----
```

The above command will generate a knative build specification as shown below:

```

apiVersion: build.knative.dev/v1alpha1
kind: BuildTemplate
metadata:
  name: build-java-maven
spec:
  parameters: ①
  - description: |
      The name of the image to push.
    name: IMAGE
  - description: |
      The context directory from where to run the build.
    name: CONTEXT_DIR
  steps:
  - args:
    - clean
    - package
    - -Duser.home=/builder/home
    - -Dimage=${IMAGE}
    image: gcr.io/cloud-builders/mvn
    name: build-maven
    volumeMounts:
    - mountPath: /builder/home/.m2
      name: m2-cache
    - mountPath: /cache
      name: kaniko-cache
    workingDir: /workspace/${CONTEXT_DIR}
  - args:
    - --context=/workspace/${CONTEXT_DIR}
    - --dockerfile=/workspace/${CONTEXT_DIR}/Dockerfile
    - --destination=${IMAGE}
    env:
    - name: DOCKER_CONFIG
      value: /builder/home/.docker
    image: gcr.io/kaniko-project/executor
    name: docker-push
    volumeMounts:
    - mountPath: /cache
      name: kaniko-cache
    workingDir: /workspace/${CONTEXT_DIR}
  volumes:
  - name: m2-cache
    persistentVolumeClaim:
      claimName: m2-cache
  - name: kaniko-cache
    persistentVolumeClaim:
      claimName: kaniko-cache

```

### maven-build-template.yaml

① The build template does not have any hard-coded values, \*the dynamic values could be passed

using **parameters** of the template.

The build done by the build steps of the template is exactly same like what was done by [raw build](#), but templates gives flexibility and ease of configuration via parameters. The template step(s) modifications are transparent to the consumers (the builds inheriting the template), the builds using the template will automatically inherit the modified steps.

## Generate Knative service

Run the following command to create the knative service that will use the image created by the **docker-build**.

```
jsonnet --ext-str contextDir='<the context dir within the application source>' \  
  --ext-str image='<the container image to use>' \  
  service-with-build-template.jsonnet \  
  | yq r - | tee ../service-with-build-template.yaml
```

(e.g.)

```
jsonnet --ext-str contextDir='java/springboot' \  
  --ext-str image='docker.io/demo/event-greeter:0.0.2' \  
  service-with-build-template.jsonnet \  
  | yq r - | tee ../service-with-build-template.yaml
```

Running the above command will result in the following knative service:

```

apiVersion: serving.knative.dev/v1alpha1
kind: Service
metadata:
  name: event-greeter
spec:
  runLatest:
    configuration:
      build:
        apiVersion: build.knative.dev/v1alpha1
        kind: Build
        spec:
          serviceAccountName: build-bot
          source:
            git:
              revision: v0.0.2 ①
              url: https://github.com/redhat-developer-demos/knative-tutorial-event-
greeter.git
          template:
            arguments: ②
            - name: IMAGE
              value: docker.io/demo/event-greeter:0.0.2
            - name: CONTEXT_DIR
              value: java/springboot
            name: build-java-maven ③
          timeout: 20m
        revisionTemplate:
          metadata:
            labels:
              app: event-greeter
          spec:
            container:
              image: docker.io/kameshsampath/event-greeter:0.0.2

```

### service-with-build-template.yaml

- ① Using another branch from the sources, just to differentiate the image produced by this build
- ② Passing the parameters required by the build template, not passing the required parameters will result in build validation failure
- ③ The name of the build template to use

## Apply resources

As we have already [applied the pre-req](#) resources, lets jump straight to applying the template and service.

## Create build template

```
oc apply -n knativetutorial -f maven-build-template.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f maven-build-template.yaml
```

## See what you have deployed

### build templates

```
# get a knative build template called build-java-maven
oc -n knativetutorial get buildtemplates.build.knative.dev build-java-maven
```

(OR)

```
kubectl -n knativetutorial get buildtemplates.build.knative.dev build-java-maven
```

### Deploy service using Build Template

With build templates the automatic build may not be triggered unless and until some service or build referring to the template gets created.

Run the following commands to deploy the service:

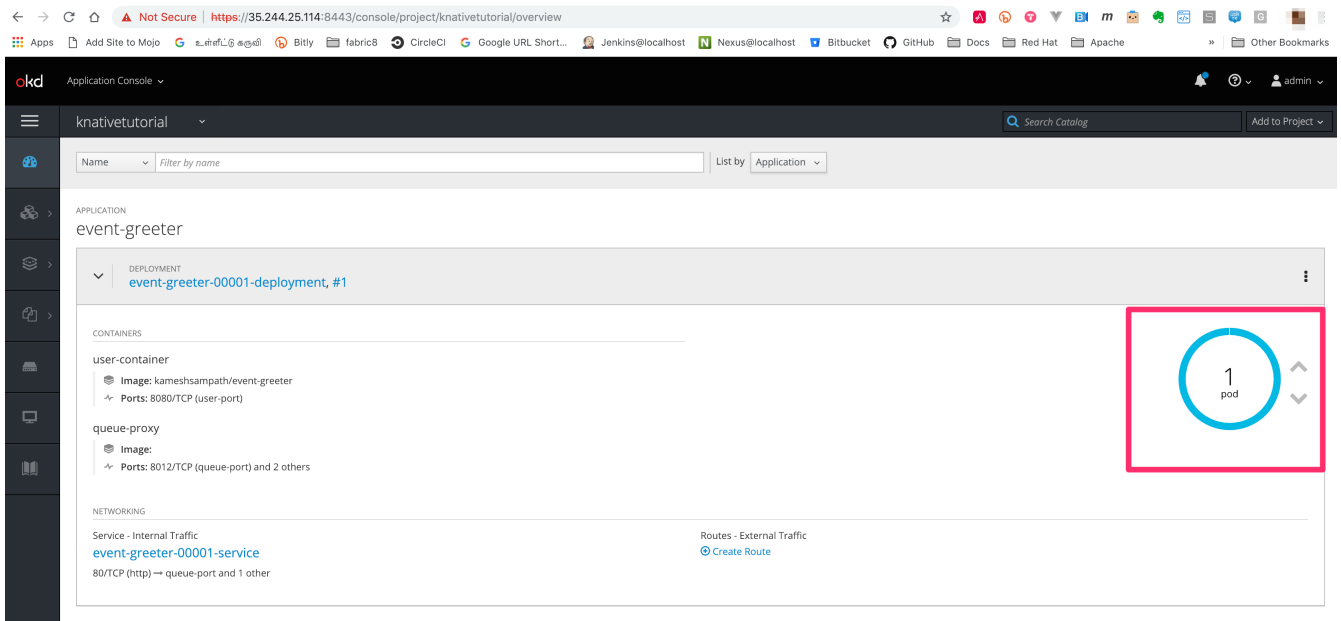
```
oc apply -n knativetutorial -f service-with-build-template.yaml
```

(OR)

```
kubectl apply -n knativetutorial -f service-with-build-template.yaml
```

You can also [watch the logs](#) of the build triggered with the service deployment.

After successful deployment of the service we should see a kubernetes deployment called **event-greeter-00001-deployment** available.



## Invoke Service

```
INGRESSGATEWAY=istio-ingressgateway
IP_ADDRESS="$(minishift ip):$(kubectl get svc $INGRESSGATEWAY --namespace istio-system
--output 'jsonpath={.spec.ports[?(@.port==80)].nodePort}')
```

```
curl -H "Host: event-greeter.knativetutorial.example.com" $IP_ADDRESS \
-H "Content-Type: application/json" \
-d '{"message": "test message"}'
```

(OR)

```
http POST $IP_ADDRESS \
'Host: event-greeter.knativetutorial.example.com' \
message="test message"
```

The above command should return an response like:

```
{
  "host": "Event greeter v0.0.2 => 'event-greeter-00001-deployment-67f96d589b-fgcj2'
: 2 ",
  "message": "test message"
}
```

## Cleanup



```
oc -n knativetutorial delete services.serving.knative.dev event-greeter
oc -n knativetutorial delete buildtemplates.build.knative.dev build-java-maven
```

**(OR)**

```
kubectl -n knativetutorial delete services.serving.knative.dev event-greeter
kubectl -n knativetutorial delete buildtemplates.build.knative.dev build-java-maven
```



You can also delete [pre-req resources](#) that were created if you dont need them any more.

# FAQs

## How to access the Knative services ?

When looking up the IP address to use for accessing your app, you need to look up the NodePort for the **istio-ingressgateway** well as the IP address used for minishift. You can use the following command to look up the value to use for the {IP\_ADDRESS} placeholder used in the sample

```
#!/bin/bash
# In Knative 0.2.x and prior versions, the `knative-ingressgateway` service was used
instead of `istio-ingressgateway`.
INGRESSGATEWAY=knative-ingressgateway

# The use of `knative-ingressgateway` is deprecated in Knative v0.3.x.
# Use `istio-ingressgateway` instead, since `knative-ingressgateway`
# will be removed in Knative v0.4.
if kubectl get configmap config-istio -n knative-serving &> /dev/null; then
    INGRESSGATEWAY=istio-ingressgateway
fi

INGRESSGATEWAY=istio-ingressgateway
IP_ADDRESS="$(minishift ip):$(kubectl get svc $INGRESSGATEWAY --namespace istio-system
--output 'jsonpath={.spec.ports[?(@.port==80)].nodePort}')"

# calling a knative service named greeter
curl -H "Host: greeter.myproject.example.com" $IP_ADDRESS
```

## Why **dev.local** suffixes for container images?

Docker images wit v2 or later format has content addressable identifier called digest. The digest remains unchanged as long the underlying image content remains unchanged.

**Source:** <https://docs.docker.com/engine/reference/commandline/images/#list-image-digests>

Let say you have a deployment like the following (note that resource definition have been trimmed for brevity):

```

apiVersion: apps/v1
kind: Deployment
metadata:
  name: helloworld
spec:
  template:
    spec:
      containers:
        - name: my-container
          image: gcr.io/knative-samples/helloworld-go
. . .

```

When you deploy this application in kubernetes, the deployment will look like:

```

apiVersion: apps/v1
kind: Deployment
metadata:
  name: helloworld
spec:
  template:
    spec:
      containers:
        - name: my-container
          image: >-
              gcr.io/knative-samples/helloworld-
go@sha256:98af362ceca8191277206b3b3854220ce125924b28a1166126296982e33882d0
. . .

```

In the above example the container image name of the deployment i.e. `gcr.io/knative-samples/helloworld-go` was resolved to its digest `gcr.io/knative-samples/helloworld-go@sha256:98af362ceca8191277206b3b3854220ce125924b28a1166126296982e33882d0`. Kubernetes extracts the digest from the image manifest. This process of resolving image tag to its digest is informally called as “Tag to Digest”.

Knative Serving deployments by default resolve the container images to digest during the deployment process. Knative Serving has been configured to skip the resolution of image name/tag to image digest for registries `ko.local` and `dev.local`, which can be used for local development builds, as the underlying image content are subject to changes during the development process.

You can also add any of the other image repo suffixes from skipping to digest by editing the configmap **config-controller** of **knative-serving** namespace and appending your repo to **registriesSkippingTagResolving** attribute of the configmap.

The following command shows an example how to add a registry **my.docker.registry** to the be skipped:



```
val=$(oc -n knative-serving get cm config-controller -oyaml | yq r -  
data.registriesSkippingTagResolving | awk '{print  
$1",my.docker.registry"}')
```

```
oc -n knative-serving get cm config-controller -oyaml | yq w -  
data.registriesSkippingTagResolving $val | oc apply -f -
```

## What is revision simpler terms?

In Knative serving, things are driven via a [Configuration](#) to make a separation between code (container images) and config. One of the good practices in configuration management is that we should always be able to rollback the application state to any “last known good configuration”, to be able to allow this type of rollback Knative creates an unique revision for each and every configuration change.