

# Scaling and Updating Applications

## Objectives

In this lab, you will:

- Scale an application with a ReplicaSet
- Apply rolling updates to an application
- Use a ConfigMap to store application configuration
- Autoscale the application using Horizontal Pod Autoscaler

## Verify the environment and command line tools

1. If a terminal is not already open, open a terminal window by using the menu in the editor: Terminal > New Terminal.

**NOTE:** It might take sometime for the Terminal Prompt to appear. In case you are unable to see the terminal prompt even after 5 minutes, please close the browser tab and relaunch the lab again.

2. Change to your project folder.

**NOTE:** If you are already in the /home/project please skip this step.

```
cd /home/project
```

3. Clone the git repository that contains the artifacts needed for this lab, if it doesn't already exist.

```
[ ! -d 'CC201' ] && https://github.com/ibm-developer-skills-network/abahi-CC_201_labs.git CC201
```

4. Change to the directory for this lab.

```
cd CC201/labs/3_K8sScaleAndUpdate/
```

5. List the contents of this directory to see the artifacts for this lab.

```
ls
```

## Build and push application image to IBM Cloud Container Registry

1. Export your namespace as an environment variable so that it can be used in subsequent commands.

```
export MY_NAMESPACE=sn-labs-$USERNAME
```

2. Use the Explorer to view the Dockerfile that will be used to build an image.

3. Build and push the image again, as it may have been deleted automatically since you completed the first lab.

```
docker build -t us.icr.io/$MY_NAMESPACE/hello-world:1 . && docker push us.icr.io/$MY_NAMESPACE/hello-world:1
```

**NOTE:** If you have tried this lab earlier, there might be a possibility that the previous session is still persistent. In such case, you will see a ‘**Layer already Exists**’ message instead of the ‘**Pushed**’ message in the above output. We would recommend you to continue with the further steps of the lab.

## Deploy the application to Kubernetes

1. Use the Explorer to edit `deployment.yaml` in this directory. The path to this file is `CC201/labs/3_K8sScaleAndUpdate/`. You need to insert your namespace where it says `<my_namespace>`. Make sure to save the file when you’re done.

**NOTE:** To know your namespace, run `echo $MY_NAMESPACE` in the terminal

2. Run your image as a Deployment.

```
kubectl apply -f deployment.yaml
```

**NOTE:** If you have tried this lab earlier, there might be a possibility that the previous session is still persistent. In such a case, you will see an ‘**Unchanged**’ message instead of the ‘**Created**’ message in the above output. We would recommend you to continue with the further steps of the lab.

3. List Pods until the status is “Running”.

```
kubectl get pods
```

**NOTE:** Please move to the next step only after you see the pod status as ‘**Running**’. In case you see ‘**Container Creating**’ as the output, please re-run the command in a few minutes.

4. In order to access the application, we have to expose it to the internet via a Kubernetes Service.

```
kubectl expose deployment/hello-world
```

5. Open a new terminal window using Terminal > New Terminal.

**NOTE:** Do not close the terminal window you were working on.

6. Cluster IPs are only accessible within the cluster. To make this externally accessible, we will create a proxy.

**Note:** This is not how you would make an application externally accessible in a production scenario.

Run this command in the new terminal window since your environment variables need to be accessible in the original window for subsequent commands.

```
kubectl proxy
```

This command will continue running until it exits. Keep it running so that you can continue to access your app.

7. Go back to your original terminal window, ping the application to get a response.

**NOTE:** Do not close the terminal window where the proxy command is still running.

```
curl -L localhost:8001/api/v1/namespaces/sn-labs-$USERNAME/services/hello-world/proxy
```

Observe the message “Hello world from hello-world-xxxxxxxxx-xxxx. Your app is up and running!

## Scaling the application using a ReplicaSet

In real-world situations, load on an application can vary over time. If our application begins experiencing heightened load, we want to scale it up to accommodate that load. There is a simple `kubectl` command for scaling.

1. Use the `scale` command to scale up your Deployment. Make sure to run this in the terminal window that is not running the proxy command.

```
kubectl scale deployment hello-world --replicas=3
```

2. Get Pods to ensure that there are now three Pods instead of just one. In addition, the status should eventually update to “Running” for all three.

```
kubectl get pods
```

3. As you did in the last lab, ping your application multiple times to ensure that Kubernetes is load-balancing across the replicas.

```
for i in `seq 10`; do curl -L localhost:8001/api/v1/namespaces/sn-labs-$USERNAME/services/hello-world/proxy; done
```

You should see that the queries are going to different Pods because of the effect of load-balancing.

4. Similarly, you can use the `scale` command to scale down your Deployment.

```
kubectl scale deployment hello-world --replicas=1
```

5. Check the Pods to see that two are deleted or being deleted.

```
kubectl get pods
```

6. Please wait for some time & run the same command again to ensure that only one pod exists.

```
kubectl get pods
```

## Perform rolling updates

Rolling updates are an easy way to update our application in an automated and controlled fashion. To simulate an update, let's first build a new version of our application and push it to Container Registry.

1. Use the Explorer to edit `app.js`. The path to this file is `CC201/labs/3_K8sScaleAndUpdate/`. Change the welcome message from '`Hello world from ' + hostname + '! Your app is up and running!\n'`' to '`Welcome to ' + hostname + '! Your app is up and running!\n'`'. Make sure to save the file when you're done.
2. Build and push this new version to Container Registry. Update the tag to indicate that this is a second version of this application. Make sure to use the terminal window that isn't running the `proxy` command.

**NOTE:** Do not close the terminal that is running the `proxy` command

```
docker build -t us.icr.io/$MY_NAMESPACE/hello-world:2 . && docker push us.icr.io/$MY_NAMESPACE/hello-world:2
```

3. List images in Container Registry to see all the different versions of this application that you have pushed so far.

```
ibmcloud cr images
```

Ensure that the new image shows `No Issues`, else re-run the image several times till there are no issues.

4. Update the deployment to use this version instead.

```
kubectl set image deployment/hello-world hello-world=us.icr.io/$MY_NAMESPACE/hello-world:2
```

5. Get a status of the rolling update by using the following command:

```
kubectl rollout status deployment/hello-world
```

6. You can also get the Deployment with the `wide` option to see that the new tag is used for the image.

```
kubectl get deployments -o wide
```

Look for the `IMAGES` column and ensure that the tag is 2.

7. Ping your application to ensure that the new welcome message is displayed.

```
curl -L localhost:8001/api/v1/namespaces/sn-labs-$USERNAME/services/hello-world/proxy
```

8. It's possible that a new version of an application contains a bug. In that case, Kubernetes can roll back the Deployment like this:

```
kubectl rollout undo deployment/hello-world
```

9. Get a status of the rolling update by using the following command:

```
kubectl rollout status deployment/hello-world
```

10. Get the Deployment with the `wide` option to see that the old tag is used.

```
kubectl get deployments -o wide
```

Look for the `IMAGES` column and ensure that the tag is 1.

11. Ping your application to ensure that the earlier ‘Hello World..Your app is up & running!’ message is displayed.

```
curl -L localhost:8001/api/v1/namespaces/sn-labs-$USERNAME/services/hello-world/proxy
```

## Using a ConfigMap to store configuration

ConfigMaps and Secrets are used to store configuration information separate from the code so that nothing is hardcoded. It also lets the application pick up configuration changes without needing to be redeployed. To demonstrate this, we’ll store the application’s message in a ConfigMap so that the message can be updated simply by updating the ConfigMap.

1. Create a ConfigMap that contains a new message.

```
kubectl create configmap app-config --from-literal=MESSAGE="This message came from a ConfigMap!"
```

**NOTE:** If you have tried this lab earlier, there might be a possibility that the previous session is still persistent. In such a case, you will see an ‘**error: failed to create configmap: configmaps “app-config” already exists**’ message, instead of the ‘**Created**’ message as below. We would recommend you to continue with the further steps of the lab.

2. Use the Explorer to edit `deployment-configmap-env-var.yaml`. The path to this file is `CC201/labs/3_K8sScaleAndUpdate/`. You need to insert your namespace where it says `<my_namespace>`. Make sure to save the file when you’re done.
3. In the same file, notice the section reproduced below. The bottom portion indicates that environment variables should be defined in the container from the data in a ConfigMap named `app-config`.

```
containers:  
- name: hello-world  
  image: us.icr.io/<my_namespace>/hello-world:3  
  ports:  
  - containerPort: 8080  
  envFrom:  
  - configMapRef:  
    name: app-config
```

4. Use the Explorer to open the `app.js` file. The path to this file is `CC201/labs/3_K8sScaleAndUpdate/`. Find the line that says, `res.send('Welcome to ' + hostname + '! Your app is up and running!\n')`.

Edit this line to look like the following:

```
res.send(process.env.MESSAGE + '\n')
```

Make sure to save the file when you're done. This change indicates that requests to the app will return the environment variable MESSAGE.

5. Build and push a new image that contains your new application code.

```
docker build -t us.icr.io/$MY_NAMESPACE/hello-world:3 . && docker push us.icr.io/$MY_NAMESPACE/hello-world:3
```

The deployment-configmap-env-var.yaml file is already configured to use the tag 3.

6. Apply the new Deployment configuration.

```
kubectl apply -f deployment-configmap-env-var.yaml
```

7. Ping your application again to see if the message from the environment variable is returned.

**NOTE:** You can run this command again. As it may not show the "This message came from a ConfigMap!" message right away.

```
curl -L localhost:8001/api/v1/namespaces/sn-labs-$USERNAME/services/hello-world/proxy
```

If you see the message, "This message came from a ConfigMap!", then great job!

**NOTE:** If your previous session is still persisting, you might see the below output. If so, we would recommend you to move to the further steps of the lab.

8. Because the configuration is separate from the code, the message can be changed without rebuilding the image. Using the following command, delete the old ConfigMap and create a new one with the same name but a different message.

```
kubectl delete configmap app-config && kubectl create configmap app-config --from-literal=MESSAGE="This message is different, and you didn't hav
```

9. Restart the Deployment so that the containers restart. This is necessary since the environment variables are set at start time.

```
kubectl rollout restart deployment hello-world
```

10. Ping your application again to see if the new message from the environment variable is returned.

```
curl -L localhost:8001/api/v1/namespaces/sn-labs-$USERNAME/services/hello-world/proxy
```

## Autoscale the hello-world application using Horizontal Pod Autoscaler

1. Replace lines 20 to 26 in the `deployment.yaml` file with the section below, under `template.spec.containers` to increase the CPU resource allocation.

```
  name: http
  resources:
    limits:
      cpu: 50m
    requests:
      cpu: 20m
```

**Note:** After making the changes, do not forget to save the file.

The updated file will be as below:

2. Apply the deployment:

```
kubectl apply -f deployment.yaml
```

3. Autoscale the `hello-world` deployment using the below command:

```
kubectl autoscale deployment hello-world --cpu-percent=5 --min=1 --max=10
```

4. You can check the current status of the newly-made `HorizontalPodAutoscaler`, by running:

```
kubectl get hpa hello-world
```

5. Please ensure that the `kubernetes proxy` is still running in the 2nd terminal. If it is not, please start it again by running:

```
kubectl proxy
```

6. Open another new terminal and enter the below command to spam the app with multiple requests for increasing the load:

```
for i in `seq 100000`; do curl -L localhost:8001/api/v1/namespaces/sn-labs-$USERNAME/services/hello-world/proxy; done
```

#### Continue further commands in the 1st terminal

7. Run the below command to observe the replicas increase in accordance with the autoscaling:

```
kubectl get hpa hello-world --watch
```

You will see an increase in the number of replicas which shows that your application has been autoscaled.

Stop this command by pressing **CTRL + C**.

8. Run the below command to observe the details of the horizontal pod autoscaler:

```
kubectl get hpa hello-world
```

You will notice that the number of replicas has increased now.

9. Stop the proxy and the load generation commands running in the other 2 terminal by pressing **CTRL + C**.

10. Delete the Deployment.

```
kubectl delete deployment hello-world
```

11. Delete the Service.

```
kubectl delete service hello-world
```

Congratulations! You have completed the lab for the third module of this course.

**© IBM Corporation. All rights reserved.**