

# Module: Cloud Computing - CSC8110

Yash Gadodia – 240679307

## Table of Contents

<b>Aim .....</b>	<b>2</b>
<b>Task 1: Deploy and access the Kubernetes Dashboard and Web Application Component ..</b>	<b>2</b>
<b>Task 2: Deploy the Monitoring Stack of Kubernetes.....</b>	<b>7</b>
<b>Task 3: Load Generator .....</b>	<b>9</b>
<b>Task 4: Monitor Benchmarking Results .....</b>	<b>13</b>
<b>Task 5: Conclusion .....</b>	<b>17</b>
<b>References.....</b>	<b>17</b>

**Aim:**

This coursework focuses on developing a comprehensive understanding and practical expertise in configuring, deploying, and managing Kubernetes-based application hosting environments. It integrates theoretical knowledge with practical skills to prepare for real-world scenarios in cloud virtualization technologies and containerized application management. Through the successful completion of Tasks 1-4 and their subsequent demonstration, hands-on experience will be gained in:

- Configuring Kubernetes environments for application hosting.
- Building, pushing, and pulling Docker images from Docker Hub.
- Creating and deploying a multi-component web application stack on Kubernetes.
- Monitoring application stack performance using Kubernetes-native tools.

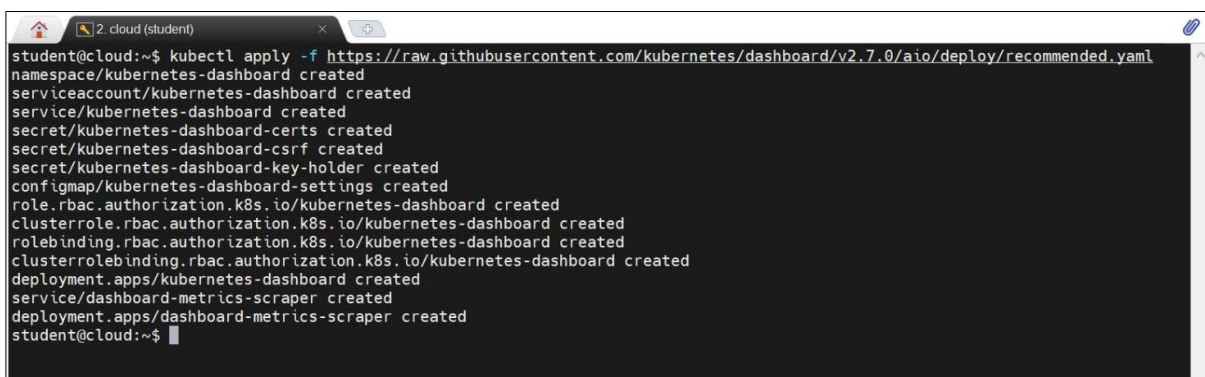
**Implementation:****Task 1: Deploy and access the Kubernetes Dashboard and a Web Application Component**

**Task Objective:** Understand and learn basic concepts of Kubernetes, and how to deploy applications to a Kubernetes cluster.

A sample *Java web application component* image, named "nclcloudcomputing/javabenchmarkapp", has been uploaded to the Docker Hub. The image contains a ready-to-use implementation of a web application deployed in a Tomcat server (an open-source web server). In terms of computational logic, the web application implements a prime number check on a large number. By doing so, the application can generate high CPU and memory load.

- a) Deploy 'Kubernetes Dashboard' on the provided VM with CLI and access/login the Dashboard.

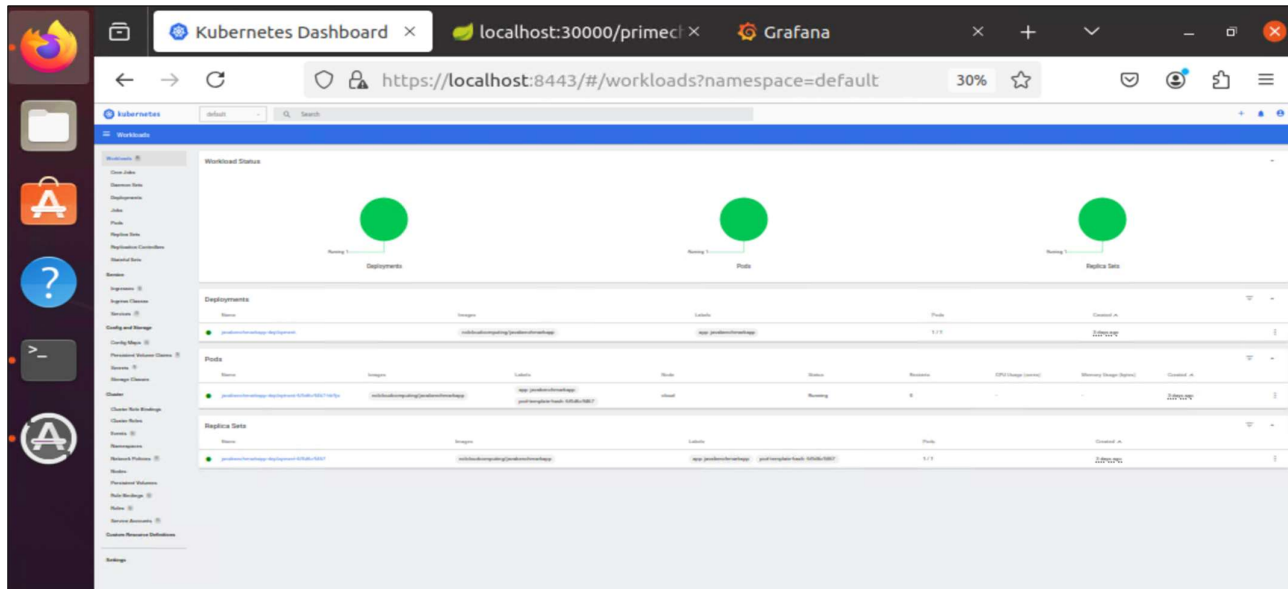
We use the following command to deploy the Kubernetes Dashboard.



```
student@cloud:~$ kubectl apply -f https://raw.githubusercontent.com/kubernetes/dashboard/v2.7.0/aio/deploy/recommended.yaml
namespace/kubernetes-dashboard created
serviceaccount/kubernetes-dashboard created
service/kubernetes-dashboard created
secret/kubernetes-dashboard-certs created
secret/kubernetes-dashboard-csrf created
secret/kubernetes-dashboard-key-holder created
configmap/kubernetes-dashboard-settings created
role.rbac.authorization.k8s.io/kubernetes-dashboard created
clusterrole.rbac.authorization.k8s.io/kubernetes-dashboard created
rolebinding.rbac.authorization.k8s.io/kubernetes-dashboard created
clusterrolebinding.rbac.authorization.k8s.io/kubernetes-dashboard created
deployment.apps/kubernetes-dashboard created
service/dashboard-metrics-scraper created
deployment.apps/dashboard-metrics-scraper created
student@cloud:~$
```



The Kubernetes Dashboard interface is displayed below.



- b) Deploy an instance of the Docker image "nclcloudcomputing/javabenchmarkapp" via CLI.

First, create a Kubernetes Deployment configuration in a YAML file to define how to deploy the Docker image and set up the required Kubernetes resources. To deploy the container from the nclcloudcomputing/javabenchmarkapp image in the Kubernetes cluster, run the following command:  
`kubectl apply -f javabenchmark_deployment.yaml`

```

1 apiVersion: apps/v1
2 kind: Deployment
3 metadata:
4   name: javabenchmarkapp-deployment
5   labels:
6     app: javabenchmarkapp
7 spec:
8   replicas: 1
9   selector:
10    matchLabels:
11      app: javabenchmarkapp
12 template:
13   metadata:
14     labels:
15       app: javabenchmarkapp
16   spec:
17     containers:
18     - name: javabenchmarkapp
19       image: nclcloudcomputing/javabenchmarkapp
20       ports:
21       - containerPort: 8080
22

```

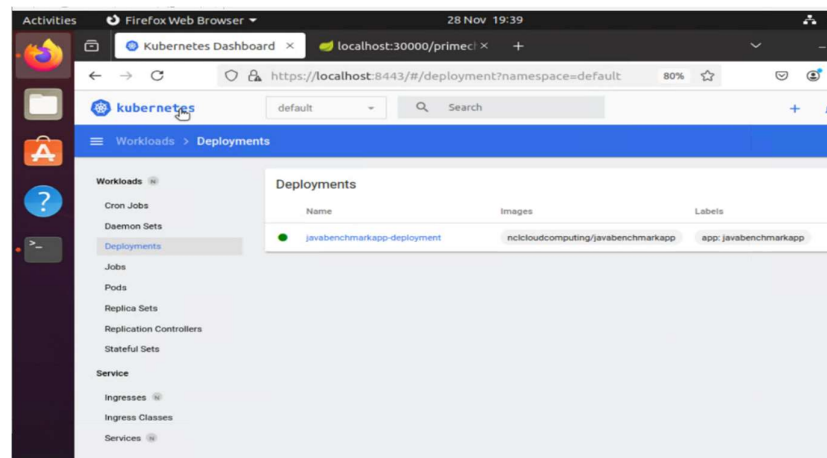
Next, define a Kubernetes Service in a YAML file to expose the application to external traffic using a NodePort. To make the application accessible externally, run: `kubectl apply -f javabenchmark_service.yaml`

```

1 apiVersion: v1
2 kind: Service
3 metadata:
4   name: javabenchmarkapp-service
5   labels:
6     app: javabenchmarkapp
7 spec:
8   type: NodePort
9   selector:
10    app: javabenchmarkapp
11   ports:
12   - protocol: TCP
13     port: 8080
14     targetPort: 8080
15     nodePort: 30000
16

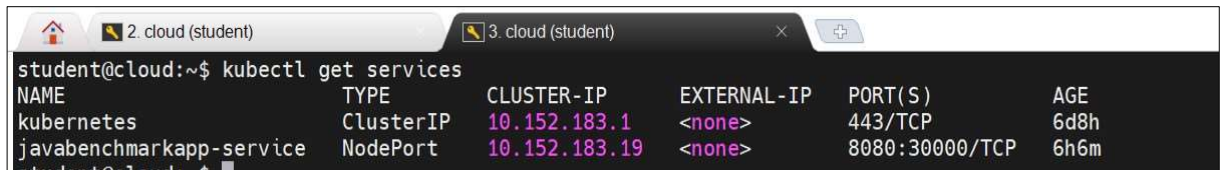
```

The command `kubectl -n kubernetes-dashboard port-forward svc/kubernetes-dashboard-kong-proxy 8443:443` is used to forward traffic from a local port (8443) to the Kubernetes service port (443), enabling access to the Kubernetes Dashboard



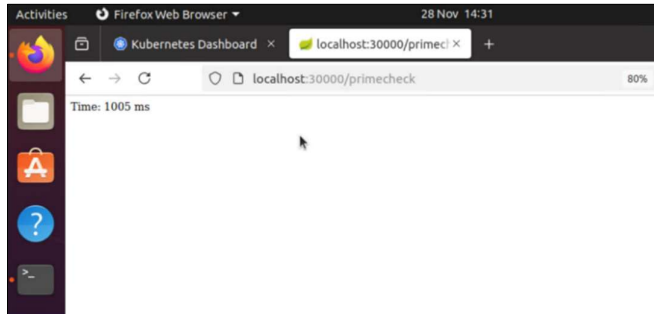
- c) Deploy a NodePort service so that the web app is accessible via <http://localhost:30000/primecheck>.

A service named `javabenchmarkapp-service` was created using the command provided below. In the configuration file, the service type is set to NodePort to make the application accessible outside the cluster.



```
student@cloud:~$ kubectl get services
```

NAME	TYPE	CLUSTER-IP	EXTERNAL-IP	PORT(S)	AGE
kubernetes	ClusterIP	10.152.183.1	<none>	443/TCP	6d8h
javabenchmarkapp-service	NodePort	10.152.183.19	<none>	8080:30000/TCP	6h6m



d) You cannot use the dashboard addon from microk8s.



## Task 2: Deploy the monitoring stack of Kubernetes

**Task Objective:** Understand and learn how to deploy a monitoring stack of Kubernetes consisting of Prometheus, metrics server, Grafana.

MicroK8s provides an easy way to enable observability features, such as Prometheus, Grafana, and the metrics server. To start with, I enabled the necessary MicroK8s addons for Prometheus and Grafana using the following command: `microk8s enable observability`

```
student@cloud:~$ microk8s enable observability
Infer repository core for addon observability
Addon core/dns is already enabled
Addon core/helm3 is already enabled
Addon core/hostpath-storage is already enabled
Enabling observability
Release "kube-prom-stack" does not exist. Installing it now.
NAME: kube-prom-stack
LAST DEPLOYED: Fri Nov 29 13:34:37 2024
NAMESPACE: observability
STATUS: deployed
REVISION: 1
NOTES:
kube-prometheus-stack has been installed. Check its status by running:
  kubectl --namespace observability get pods -l "release=kube-prom-stack"

Visit https://github.com/prometheus-operator/kube-prometheus for instructions on how to create & configure Alertmanager and
Prometheus instances using the Operator.
Release "loki" does not exist. Installing it now.
NAME: loki
LAST DEPLOYED: Fri Nov 29 13:35:09 2024
NAMESPACE: observability
STATUS: deployed
REVISION: 1
NOTES:
The Loki stack has been deployed to your cluster. Loki can now be added as a datasource in Grafana.

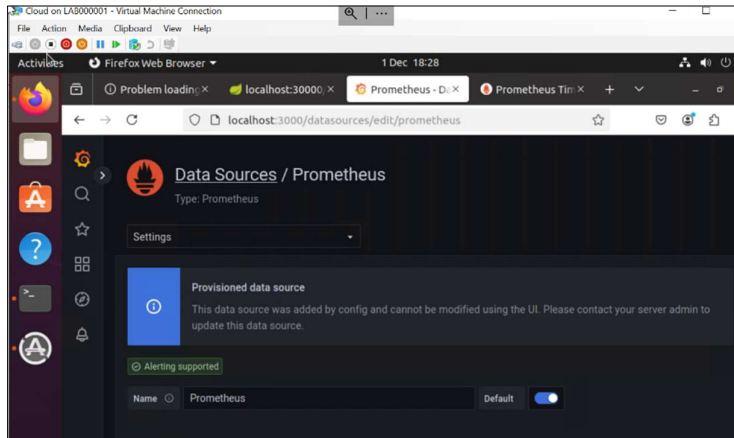
See http://docs.grafana.org/features/datasources/loki/ for more detail.
Release "tempo" does not exist. Installing it now.
NAME: tempo
LAST DEPLOYED: Fri Nov 29 13:35:15 2024
NAMESPACE: observability
STATUS: deployed
```

```
student@cloud:~$ microk8s kubectl get pods -n observability
```

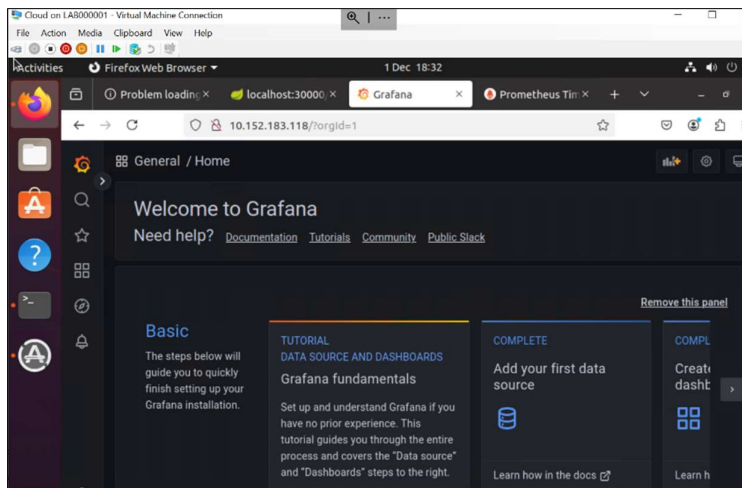
NAME	READY	STATUS	RESTARTS	AGE
kube-prom-stack-kube-prometheus-operator-64ffd55b77-kvgvgv	1/1	Running	0	6m12s
kube-prom-stack-prometheus-node-exporter-96lpd	1/1	Running	0	6m13s
kube-prom-stack-kube-state-metrics-6c586bf4c8-nvjhz	1/1	Running	0	6m12s
kube-prom-stack-grafana-6c47f548d6-kj5k4	3/3	Running	0	6m13s
alertmanager-kube-prom-stack-kube-prometheus-alertmanager-0	2/2	Running	1 (5m7s ago)	5m53s
loki-promtail-25cv4	1/1	Running	0	5m54s
tempo-0	2/2	Running	0	5m49s
prometheus-kube-prom-stack-kube-prometheus-prometheus-0	2/2	Running	0	5m53s
loki-0	1/1	Running	0	5m54s

**Explanation of Task2:** By default, the Grafana and Prometheus services are not directly accessible from the host system. To make these services reachable, I used the `kubectl port-forward` command to expose them to specific ports on the local machine. For Prometheus, I forwarded port 9090 by running the command: `microk8s kubectl port-forward -n observability service/kube-prom-stack-kube-prometheus 9090:9090`. Similarly, for Grafana, I forwarded port 3000 with the following command: `microk8s kubectl port-forward -n observability service/kube-prom-stack-grafana 3000:80`. I attach the snip of all the outputs below.

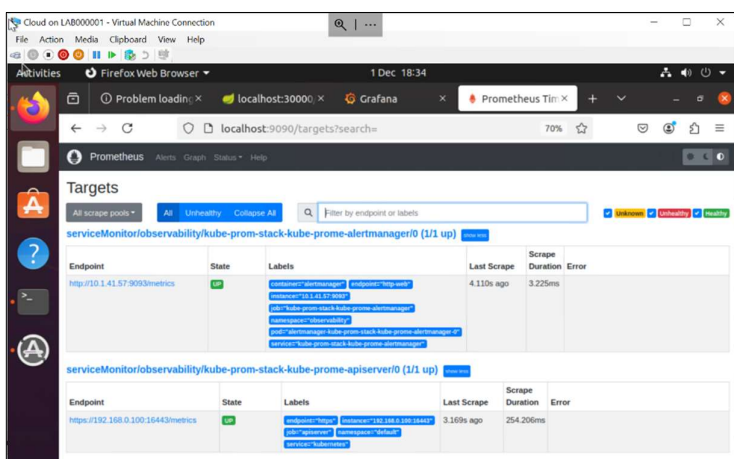
## a) Grafana must reach Prometheus as a data source



## b) Grafana must be reachable from the host



## c) Prometheus must reach each instance of the metric server to collect metrics



## d) The metric server must be deployed as a Daemon Set

```
student@cloud:~$ microk8s kubectl get deployments,daemonsets -A | grep -E "kube-state-metrics|node-exporter"
observability deployment.apps/kube-prom-stack-kube-state-metrics 1/1 1 1 1 1 2d5h
observability daemonset.apps/kube-prom-stack-prometheus-node-exporter 1 1 1 1 1 1 2d5h
```



### Task 3: Load Generator

**Task Objective:** Understand the logic of the load generator of benchmarking web applications and the process of deploying your own application of the cluster. Additionally, understand how to build and push a Docker image from scratch.

- a) Write a load generator with the following specifications
  - a. Accepts two configurable values either via a config file or environment variables. target (The address for the load generation) and frequency (Request per second)
  - b. Generate web request to the target at the specified frequency
  - c. Collect 2 types of metrics. Average response time and accumulated number of failures
  - d. Request should timeout if it takes more than 10 seconds. Counted as failures
  - e. Test results need to be printed to the console
  - f. There are no requirements in programming language

```
* load_generator.py
1 import requests
2 import time
3 import os
4
5 class LoadGenerator:
6     def __init__(self, target=None, frequency=None):
7         self.target = target or os.environ.get("target", "http://192.168.0.100:30000/primecheck")
8         self.frequency = frequency or float(os.environ.get("Frequency", 10.0))
9         self.failures = 0
10        self.total_response_time = 0
11        self.total_requests = 0
12
13    def perform_request(self): """HTTP GET request to the target."""
14        try:
15            response = requests.get(self.target, timeout=10)
16            response.raise_for_status()
17            return response
18        except requests.exceptions.RequestException as e:
19            print(f"Request failed: {e}")
20            self.failures += 1
21            return None
22
23    def calculate_response_time(self, start_time): """Calculate and return the response time."""
24        end_time = time.time()
25        return end_time - start_time
26
27    def update_metrics(self, response_time): """Update the metrics for total response time and total requests."""
28        self.total_response_time += response_time
29        self.total_requests += 1
30
31    def print_metrics(self): """Print the collected metrics."""
32        if self.total_requests > 0:
33            avg_response_time = self.total_response_time / self.total_requests
34            print(f"Average Response Time: {avg_response_time:.2f} seconds")
35            print(f"Total Failures: {self.failures}")
36            print(f"Total Requests: {self.total_requests}")
37
38    def generate_load(self): """Generate load by repeatedly performing requests."""
39        while True:
40            start_time = time.time()
41            response = self.perform_request()
42
43            if response:
44                response_time = self.calculate_response_time(start_time)
45                self.update_metrics(response_time)
46                print(f"Request successful. Response Time: {response_time:.2f} seconds")
47
48            time.sleep(1 / self.frequency)
49            self.print_metrics()
50
51 if __name__ == "__main__":
52     load_generator = LoadGenerator()
53     load_generator.generate_load()
54
```

This Python script functions as a load generator that sends HTTP GET requests to a specified target URL at a user-defined frequency. The target URL and request frequency can be configured through environment variables, with fallback default values if the variables are not set. The script ensures requests are sent at the desired rate by using `time.sleep(1 / frequency)`. It monitors two main metrics: the average response time, which is calculated by dividing the total response time by the number of successful requests, and the total failures, which increases each time a request times out or encounters an error. To manage timeouts, the `requests.get()` method is set with a 10-second timeout, and any request that exceeds this limit is treated as a failure. At regular intervals, the script outputs test results to the console, including the average response time, the number of failures, and the total number of requests, allowing for real-time performance monitoring of the target system during the load test.

- b) After programming, pack the program as a standalone Docker image and push it to the local registry at port 32000. Name the image as *load-generator*.

First, we start a Docker registry on our host, running on port 32000, using the **command**: `docker run -d -p 32000:5000 --name registry`.

```
student@cloud:~$ docker run -d -p 32000:5000 --name registry registry:2
Unable to find image 'registry:2' locally
2: Pulling from library/registry
dc0decf4841d: Pull complete
6cb0aa443e23: Pull complete
813676e291ef: Pull complete
dc2fb7dcec61: Pull complete
916205650bfe: Pull complete
Digest: sha256:543dade69668e02e5768d7ea2b0aa4fae6aa7384c9a5a8dbec2be5136079ddb
Status: Downloaded newer image for registry:2
e3e530ac9ecddc12e787d8e66119095f83755da257e8be3400b9fd3a96c7e249
student@cloud:~$ docker ps
```

CONTAINER ID	IMAGE	COMMAND	CREATED	STATUS	PORTS
e3e530ac9ecd	registry:2	"/entrypoint.sh /etc..."	12 seconds ago	Up 9 seconds	0.0.0.0:32000→5000/tcp, :::32000→5000/tcp

Next, we build the load-generator Docker image with the following command: `docker build -t load-generator .`

```

dockerfile
1 # Use the official Python 3.9 slim image as the base image
2 FROM python:3
3
4 # Set the working directory inside the container
5 WORKDIR /app
6
7 # Copy the Python script into the container
8 COPY load_generator.py .
9
10 # Install the required Python library
11 RUN pip install requests
12
13 EXPOSE 30000
14
15 # Set the default command to execute the load generator script
16 CMD ["python", "load_generator.py"]
17

```

```

student@cloud:~$ docker build -t load-generator .
Sending build context to Docker daemon 240.3MB
Step 1/5 : FROM python:3
--> c41ea8273365
Step 2/5 : WORKDIR /app
--> Running in 3620ff109dc6
Removing intermediate container 3620ff109dc6
--> d8bef7fa704a
Step 3/5 : COPY load_generator.py .
--> a5ab1abda056
Step 4/5 : RUN pip install requests
--> Running in dd964cd14f68
Collecting requests
  Downloading requests-2.32.3-py3-none-any.whl.metadata (4.6 kB)
Collecting charset-normalizer<4, ≥2 (from requests)
  Downloading charset-normalizer-3.4.0-cp313-cp313-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (34 kB)
Collecting idna<4, ≥2.5 (from requests)
  Downloading idna-3.10-py3-none-any.whl.metadata (10 kB)
Collecting urllib3<3, ≥1.21.1 (from requests)
  Downloading urllib3-2.2.3-py3-none-any.whl.metadata (6.5 kB)
Collecting certifi≥2017.4.17 (from requests)
  Downloading certifi-2024.8.30-py3-none-any.whl.metadata (2.2 kB)
Downloading requests-2.32.3-py3-none-any.whl (64 kB)
Downloading certifi-2024.8.30-py3-none-any.whl (167 kB)
Downloading charset-normalizer-3.4.0-cp313-cp313-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (144 kB)
Downloading idna-3.10-py3-none-any.whl (70 kB)
Downloading urllib3-2.2.3-py3-none-any.whl (126 kB)
Installing collected packages: urllib3, idna, charset-normalizer, certifi, requests
Successfully installed certifi-2024.8.30 charset-normalizer-3.4.0 idna-3.10 requests-2.32.3 urllib3-2.2.3
WARNING: Running pip as the 'root' user can result in broken permissions and conflicting behaviour with the system package manager, possibly rendering your system unusable. It is recommended to use a virtual environment instead: https://pip.pypa.io/warnings/venv. Use the --root-user-action option if you know what you are doing and want to suppress this warning.

[notice] A new release of pip is available: 24.2 → 24.3.1
[notice] To update, run: pip install --upgrade pip
Removing intermediate container dd964cd14f68
--> fbe94d1f7fec
Step 5/5 : CMD ["python", "load_generator.py"]

```

After building the image, we can see load-generator listed in the output of the `docker images` command.

```

student@cloud:~$ docker images
REPOSITORY          TAG          IMAGE ID          CREATED           SIZE
load-generator      latest      9c5568a1b457     44 seconds ago   1.03GB
python             3          c41ea8273365     6 weeks ago      1.02GB
registry           2          c18a86d35e98     14 months ago    25.4MB

```

Once the image is created, we tag it and push it to our Docker registry. To tag the image, we use the command: `docker tag load-generator localhost:32000/load-generator`. Finally, we push the image to the registry using: `docker push localhost:32000/load-generator`.

```
student@cloud:~$ docker tag load-generator localhost:32000/load-generator
student@cloud:~$ docker push localhost:32000/load-generator
Using default tag: latest
The push refers to repository [localhost:32000/load-generator]
ed2757717041: Pushed
953b03b0ea16: Pushed
20da8de59e62: Pushed
1e8b53dc33fa: Layer already exists
94878487f485: Layer already exists
20a1c75bf72e: Layer already exists
96d99c63b722: Layer already exists
00547dd240c4: Layer already exists
b6ca42156b9f: Layer already exists
24b5ce0f1e07: Layer already exists
latest: digest: sha256:7690962e920fd37a4df8367b3c450bcfa4041afa2a6b94fc6b2a2a95c9b3d530 size: 2420
```

After pushing, the `docker images` command will show both the `load-generator` and `localhost:32000/load-generator` images.

```
student@cloud:~$ docker images
```

REPOSITORY	TAG	IMAGE ID	CREATED	SIZE
load-generator	latest	9c5568a1b457	4 minutes ago	1.03GB
localhost:32000/load-generator	latest	9c5568a1b457	4 minutes ago	1.03GB
python	3	c41ea8273365	6 weeks ago	1.02GB
registry	2	c18a86d35e98	14 months ago	25.4MB



## Task 4: Monitor benchmarking results

**Task Objective:** To learn and understand how to monitor container metrics.

- a) Deploy *load-generator* service created in Task 3.

```
* load-generator.yaml
1 kind: Deployment
2 apiVersion: apps/v1
3 metadata:
4   name: load-generator
5 spec:
6   replicas: 1
7   selector:
8     matchLabels:
9       name: load-generator
10  template:
11    metadata:
12      name: testpod1
13      labels:
14        name: load-generator
15    spec:
16      containers:
17        - name: c01
18          image: localhost:32000/load-generator
19          ports:
20            - containerPort: 30000
21
22 ---
23 kind: Service
24 apiVersion: v1
25 metadata:
26   name: demoservice
27 spec:
28   ports:
29     - port: 30000
30       targetPort: 30000
31   selector:
32     name: load-generator
33   type: NodePort
34
```

```
student@cloud:~$ kubectl apply -f load-generator-deployment.yaml
deployment.apps/load-generator created
```

Above YAML file configures a Kubernetes Deployment and Service to deploy and expose a load generator application. The Deployment defines a single replica of a pod, labeled Load-generator, which runs a container named c01 with the image localhost:32000/load-generator. The container listens on port 30000. The Service, configured as a NodePort, exposes the pods on port 30000, enabling external access to the load generator. It links to the pods using the label selector name: Load-generator.

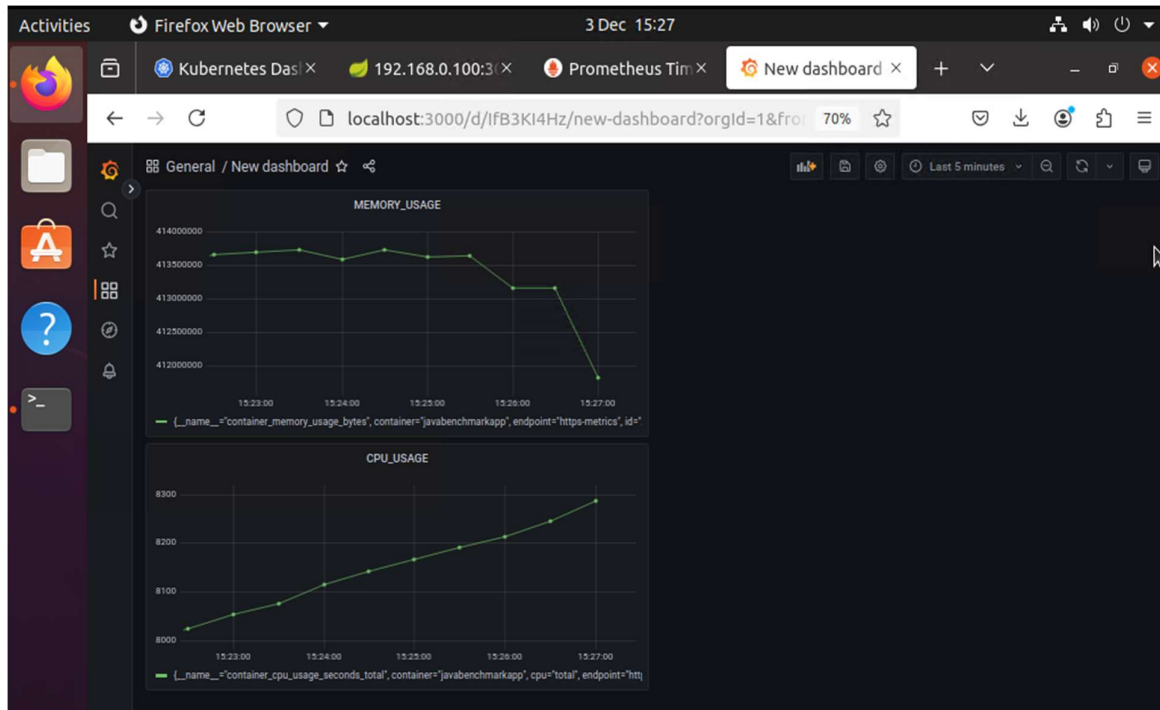
We deploy the load generator by configuring environment variable values that the load-generator code uses to repeatedly send requests to our JavaBenchmarkApp. The figure below shows the logs generated after deploying the service.

```
student@cloud:~$ kubectl logs load-generator-6556685d9d-pkpgp
Request successful. Response Time: 1.03 seconds
Average Response Time: 1.03 seconds
Total Failures: 0
Total Requests: 1
Request successful. Response Time: 0.86 seconds
Average Response Time: 0.94 seconds
Total Failures: 0
Total Requests: 2
Request successful. Response Time: 0.84 seconds
Average Response Time: 0.91 seconds
Total Failures: 0
Total Requests: 3
Request successful. Response Time: 0.92 seconds
Average Response Time: 0.91 seconds
Total Failures: 0
Total Requests: 4
Request successful. Response Time: 0.98 seconds
Average Response Time: 0.93 seconds
Total Failures: 0
Total Requests: 5
Request successful. Response Time: 0.97 seconds
Average Response Time: 0.93 seconds
Total Failures: 0
Total Requests: 6
Request successful. Response Time: 0.88 seconds
Average Response Time: 0.93 seconds
Total Failures: 0
Total Requests: 7
Request successful. Response Time: 0.86 seconds
Average Response Time: 0.92 seconds
Total Failures: 0
Total Requests: 8
Request successful. Response Time: 0.90 seconds
Average Response Time: 0.92 seconds
Total Failures: 0
Total Requests: 9
Request successful. Response Time: 0.91 seconds
```



- b) During the benchmarking, create a new dashboard on Grafana and add 2 new panels which should contain queries of CPU/memory usage of the web application

I created a new dashboard on Grafana and added two new panels displaying queries for CPU and memory usage of the web application, which you can see in the figure below.

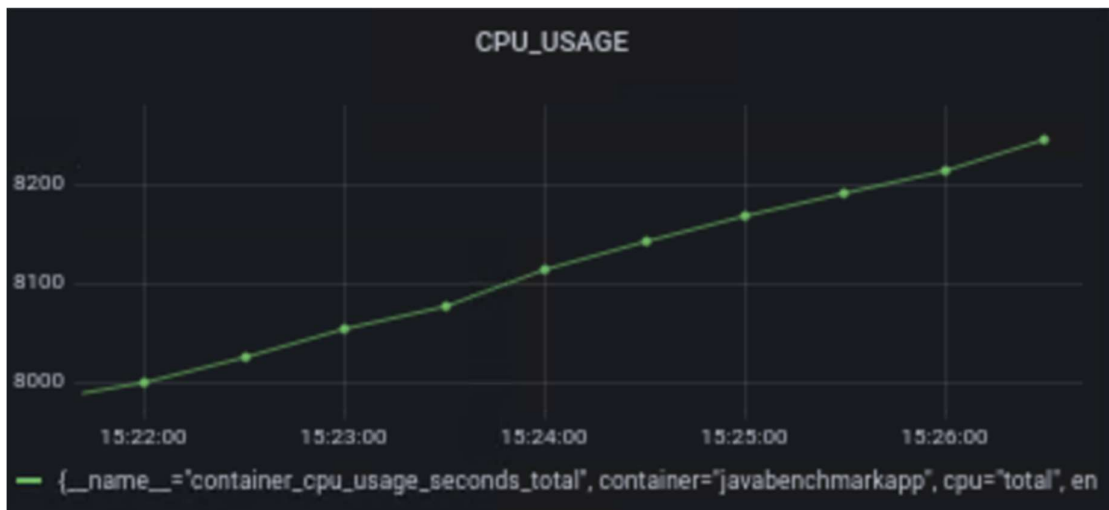


- c) Screenshot the two panels  
i) Cpu\_Usage

The image below presents the metrics for `container_cpu_usage_seconds_total`, which illustrates a steady rise in CPU usage over time. The graph reflects an increase in CPU utilization as the web application starts handling requests. CPU usage only decreases when there is a request failure or an issue with the service. However, in this case, since there are no request failures, the graph shows a consistent upward trend.

Server Metric: `container_cpu_usage_seconds_total`  
Label: `container > javabenchmarkapp`

Snapshot:



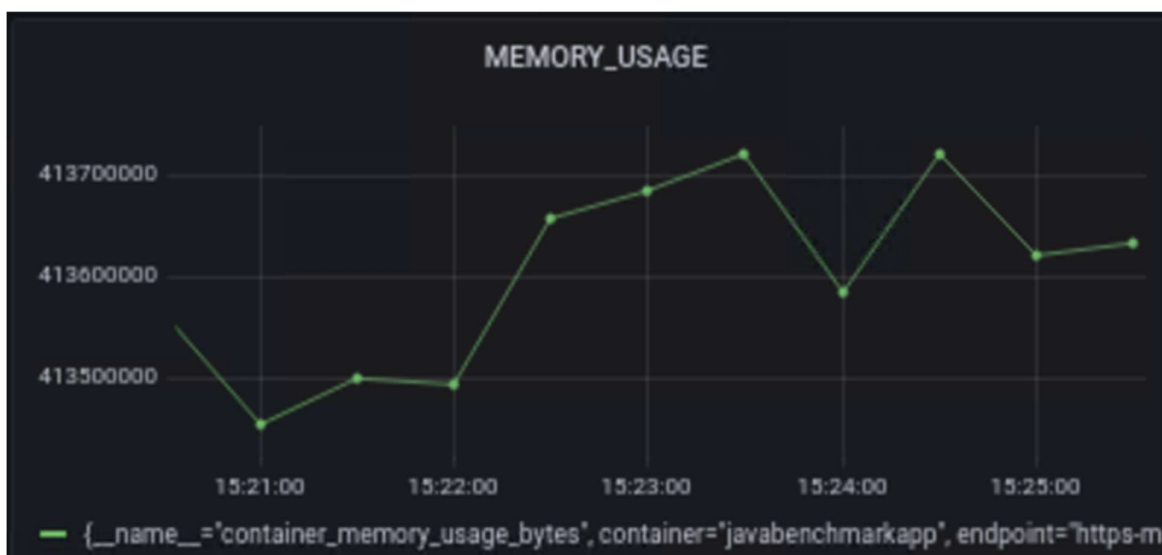
## ii) Memory\_usage

The graph below displays the `container_memory_usage_bytes` metric, illustrating a non-linear pattern of memory consumption over time. Memory usage is influenced by the behaviour of the web application. When the `javabenchmarkapp` service is down, the graph shows a decline in memory usage, while an increase in memory usage occurs when the service is operational.

Server Metric: `container_memory_usage_bytes`

Label: `container > javabenchmarkapp`

Snapshot:



**Conclusion:**

Throughout this coursework, we gained hands-on experience in deploying, managing, and monitoring applications within a Kubernetes environment. We also developed a load generation tool for benchmarking, which provided valuable insights into application performance and resource usage, helping to identify areas for optimization and scalability. By utilizing Grafana, we created a comprehensive dashboard to visualize CPU and memory usage, revealing a linear relationship between increased load and resource consumption. This exercise has enhanced our understanding of the impact of load on application performance and the importance of monitoring tools in optimizing cloud-based applications. Furthermore, it emphasized the need for effective scaling strategies to handle varying levels of demand.

**References:**

- a) [Kubernetes](#): Kubernetes Documentation and Access the Kubernetes Cluster Web UI Dashboard.
- b) [Stackoverflow](#): Stack Overflow, for various doubts related to implementation and errors