



**B.M.S. COLLEGE OF ENGINEERING**  
Autonomous Institute, Affiliated to VTU  
Estd. 1946

**DEPARTMENT OF CSE**

**CTY Project Work In collaboration with HPE**

**Weekly report**

<b>Project Title</b>	Comparative study of container runtimes				
<b>Student Team</b>	Name: Kizhakel Sharat Prasad  USN: 1BM19CS074  SEM:UG-6th sem	Name: Sharan S Pai  USN: 1BM19CS146  SEM:UG-6th sem			
	Name: Disha N  USN: 1BM19CS051  SEM:UG-6th sem	Name: Vallisha M  USN: 1BM19CS177  SEM:UG-6th sem			
<b>Faculty Mentor</b>	Dr. Pallavi G B  (Assistant Professor)  Prof. Madhavi R P  (Associate Professor)	<b>HPE Mentors</b>	Ravi Mehta		
<b>Review for the Period</b>	30-06-2022	22-07-2022			
<b>Github Link</b>	Application tested: <a href="https://github.com/sharan-s-pai/dockerized-goals-app">https://github.com/sharan-s-pai/dockerized-goals-app</a>  Reports: <a href="https://github.com/Sharat-Kizhakel/HPE-CTY">https://github.com/Sharat-Kizhakel/HPE-CTY</a>				
<b>Work done so far</b>	<ol style="list-style-type: none"><li>Successfully installed five runtimes discussed under problem statement</li><li>Plotted their cpu usages, memory as well as network usage using prometheus</li><li>exported data sources from prometheus and plotted relevant dashboards in grafana</li></ol>				

## Introduction

### **What is a container runtime?**

The container runtime is the low-level component that creates and runs containers. Docker currently uses runC, the most popular runtime, which adheres to the OCI standard that defines container image formats and execution. However, because Docker observes OCI-compliance, any OCI-compliant runtime should work. It is software that runs and manages the components required to run containers. They can be broadly classified into two main categories:

#### 1. Open Container Initiative (OCI):

The Open Container Initiative (**OCI**) is a Linux Foundation project to design open standards for containers. Established in June 2015 by Docker and other leaders in the container industry.

OCI currently contains two specifications:

- the Runtime Specification (runtime-spec)
- Image Specification (image-spec).

OCI runtime spec defines how to run the OCI image bundle as a container. OCI image spec defines how to create an OCI Image, which includes an image manifest, a filesystem (layer) serialization, and an image configuration.

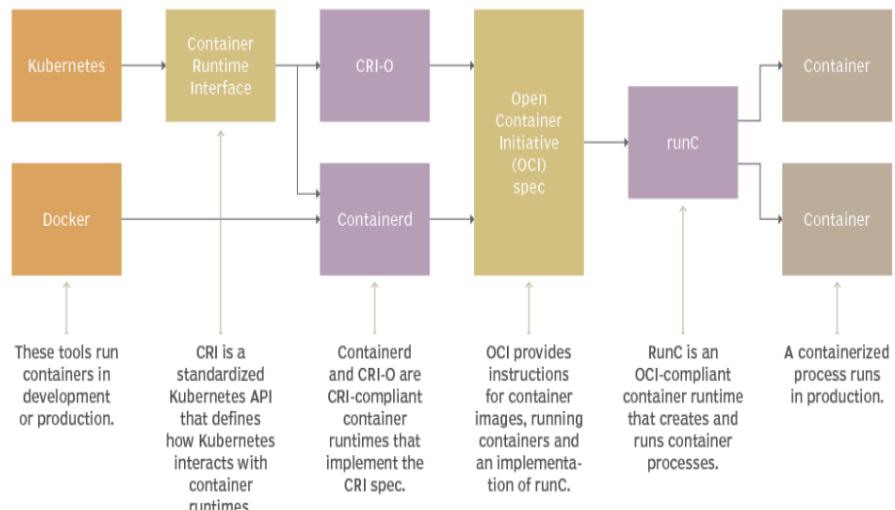
#### 2. Container Runtime Interface (CRI):

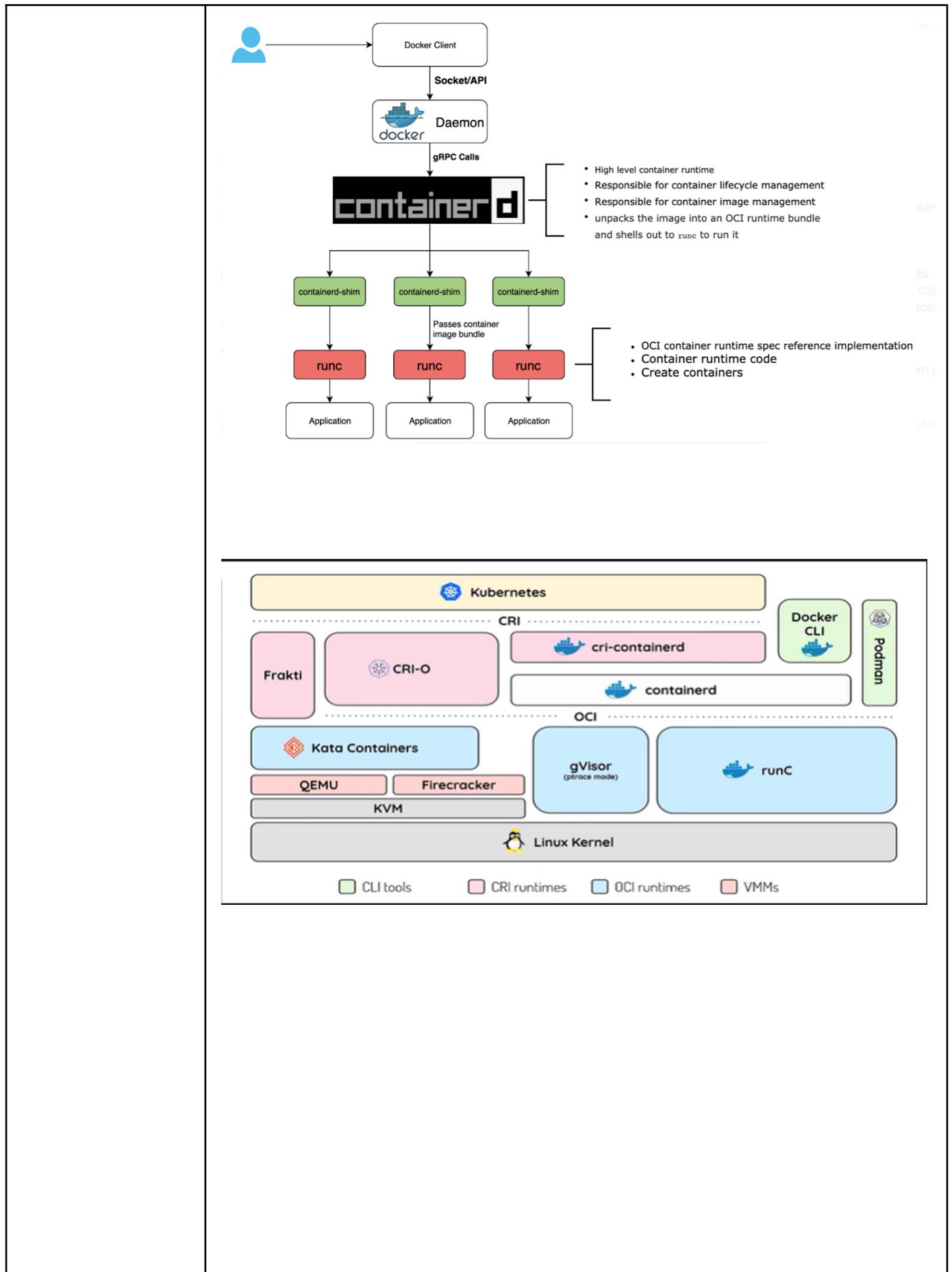
Container runtime interface (CRI) is a plugin interface that lets the kubelet—an agent that runs on every node in a Kubernetes cluster—use more than one type of container runtime. The CRI has additional concerns over an OCI Runtime including image management and distribution, storage, snapshotting, networking (distinct from the CNI), and more.

- Open Container Initiative (OCI) Runtimes
  - Native Runtimes

- runC
- Railcar
- Crun
- rkt
- Sandboxed and Virtualized Runtimes
  - gvisor
  - nabla-containers
  - runV
  - clear containers
  - kata-containers
- Container Runtime Interface (CRI)
  - containerd
  - cri-o

How Docker, Kubernetes, OCI, CRI-O, containerd and runC collaborate to run containers





## Runtimes which we have taken into consideration?

### 1. runc:

Open Container Initiative runtime runc is a command line client for running applications packaged according to the Open Container Initiative (OCI) format and is a compliant implementation of the Open Container Initiative specification. Runc is a so-called “container runtime”.

### 2. crun:

It is a lightweight fully featured OCI runtime and C library for running containers. It is faster than runc and has a much lower memory footprint.

### 3. runsc:

gVisor is an application kernel, written in Go, that implements a substantial portion of the Linux system surface. It includes an Open Container Initiative (OCI) runtime called runsc that provides an isolation boundary between the application and the host kernel. The runsc runtime integrates with Docker and Kubernetes, making it simple to run sandboxed containers.

### 4. sysbox-runc:

Sysbox is an open-source and free container runtime (a specialized “runc”), developed by Nestybox, that enhances containers in two key ways:

- Improves container isolation
- Enables containers to run same workloads as VMs

### 5. kata:

kata-runtime, referred to as “the runtime”, is the Command-Line Interface (CLI) part of the Kata Containers runtime component. It leverages the virtcontainers package to provide a high-performance standards-compliant runtime that creates hardware-virtualized Linux containers running on Linux hosts. The runtime is OCI-compatible, CRI-O-compatible, and Containerd-compatible, allowing it to work seamlessly with both Docker and Kubernetes respectively.

## Application Overview

In the course of the Project the same application has been tested with different runtimes. The application in question is a simple react app which comprises of 3 main containers which we monitored namely:

1. node-server
  2. docker
  3. mongodb

The app contains a docker-compose file with all relevant dependencies including a cAdvisor and Prometheus port numbers. In addition to this bind mounts have been used to persist data. The performance of a container cannot be tested without sending a significant payload to the server. For this purpose a movies dataset has been used. The runtime was changed everytime we editing the installed runtime in the available runtimes in the daemon.json file which comes installed by default with docker. The sure test for whether we are using the new runtime can be confirmed by checking the kernel version using the uname option in docker.

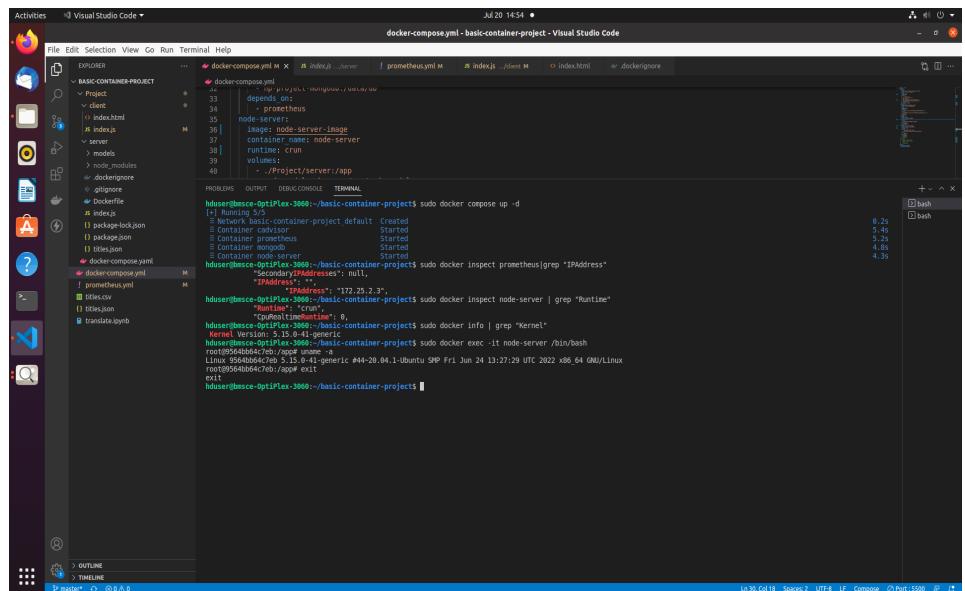
The following are the screenshots proving successful installation of runtimes:

## 1. runc

The screenshot shows a Visual Studio Code interface with the following details:

- Activity Bar:** Activities, Visual Studio Code.
- File Explorer:** Shows the project structure for "BASIC-CONTAINER-PROJECT".
- Terminal:** Shows the command `sudo docker compose up -d` being run, followed by the output of the Docker Compose up command which creates services for prometheus, advisor, and mongo.
- Output:** Shows logs from the Docker container, including kernel version and a grep command for IP addresses.
- Problems:** No problems listed.
- Search:** Placeholder "Search" and a search bar.
- Outline:** Shows the file structure.
- Timeline:** Shows a timeline of recent changes.
- Status Bar:** Master, 0 @ 0.0, exit, Lin 38, Col 18, Spaces: 2, UFT-8, LF, Compose, Port: 15600.

## 2. crun



The screenshot shows a Visual Studio Code interface with the following details:

- File Explorer:** Shows a project named "BASIC-CONTAINER-PROJECT" containing files like dockerfile, index.js, package.json, package-lock.json, prometheus.yml, Dockerfile, and docker-compose.yml.
- Terminal:** Displays the command `sudo docker compose up -d` being run, followed by the output of `sudo docker inspect` for the prometheus and node-server containers. The output includes container IDs, IP addresses (e.g., 172.25.2.3), and kernel versions (e.g., 5.15.0-41-generic).
- Status Bar:** Shows the terminal has 30 columns, 19 rows, and is using UTF-8 encoding.

## 3. runsc

Activities Visual Studio Code

File Edit Selection View Go Run Terminal Help

OPEN EDITORS

- index.js basic-container-project/pr... M
- docker-compose.yml basic-container-project M
- PROJECT
- basic-container-project
- Project > client
- server > models
- node\_modules
- .dockignore
- .gitignore
- Dockerfile
- index.js JS
- package-lock.json
- package.json
- titles.json
- translate.ipynb
- docker-compose.yml M
- prometheus.yml M
- titles.csv M
- titles.json M
- translate.ipynb M

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL

```
hduser@bmse-OptiPlex-3060:~/project/basic-container-project$ sudo docker compose up -d
[+] Running 5/5
  ● Network basic-container-project default Created
  ● Container cAdvisor Started
  ● Container prometheus Started
  ● Container mongodb Started
  ● Container node-server Started
hduser@bmse-OptiPlex-3060:~/project/basic-container-project$ sudo docker inspect mongo | grep "IPAddress"
"SecondaryIPAddresses": null,
"IPAddress": "172.24.8.4",
hduser@bmse-OptiPlex-3060:~/project/basic-container-project$ sudo docker inspect prometheus | grep "IPAddress"
"SecondaryIPAddresses": null,
"IPAddress": "",
"IPAddress": "172.24.8.3",
hduser@bmse-OptiPlex-3060:~/project/basic-container-project$ sudo docker inspect node-server | grep "Runtime"
"Runtime": "runc",
"cpuRealtimeRuntime": 0,
hduser@bmse-OptiPlex-3060:~/project/basic-container-project$ sudo docker info | grep "Kernel"
WARNING: No swap limit support
Kernel Version: 4.15.0-189-generic
hduser@bmse-OptiPlex-3060:~/project/basic-container-project$ sudo docker exec -it node-server /bin/bash
root@6dd0ce1c00c:/app# uname -a
Linux 6dd0ce1c00c 4.4.0 #1 SMP Sun Jan 10 15:06:54 PST 2016 x86_64 GNU/Linux
root@6dd0ce1c00c:/app# exit
exit
hduser@bmse-OptiPlex-3060:~/project/basic-container-project$
```

Ln 35, Col 20 Spaces:2 UFF:8 LF Compose ⚡ Port:5500 ⚡

#### 4. sysbox-runc

Activities Visual Studio Code

File Edit Selection View Go Run Terminal Help

OPEN EDITORS

- index.js basic-container-project M
- index.html M
- prometheus.yml M
- index.js - violent M
- index.html M
- dockersigns
- docker-compose.yml M
- prometheus.yml M
- titles.csv M
- titles.json M
- translate.ipynb M

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL

```
hduser@bmse-OptiPlex-3060:~/basic-container-project$ sudo docker compose down
[+] Running 5/5
  ● Container node-server Removed
  ● Container mongoDB Removed
  ● Container prometheus Removed
  ● Container cAdvisor Removed
  ● Network basic-container-project default Removed
hduser@bmse-OptiPlex-3060:~/basic-container-project$ sudo docker compose up -d
[+] Running 5/5
  ● Container mongoDB Created
  ● Container prometheus Started
  ● Container cAdvisor Started
  ● Container node-server Started
  ● Network basic-container-project default Created
hduser@bmse-OptiPlex-3060:~/basic-container-project$ sudo docker inspect prometheus | grep "IPAddress"
"SecondaryIPAddresses": null,
"IPAddress": "172.25.1.3",
hduser@bmse-OptiPlex-3060:~/basic-container-project$ sudo docker info | grep "Kernel"
Kernel Version: 5.15.0-41-generic
hduser@bmse-OptiPlex-3060:~/basic-container-project$ sudo docker exec -it node-server /bin/bash
root@24311dc63079:/app# uname -a
Linux 24311dc63079 5.15.0-41-generic #44~20.04.1-Ubuntu SMP Fri Jun 24 13:27:29 UTC 2022 x86_64 GNU/Linux
root@24311dc63079:/app# exit
exit
hduser@bmse-OptiPlex-3060:~/basic-container-project$ sudo docker inspect node-server | grep "Runtime"
"Runtime": "sysbox-runc",
"cpuRealtimeRuntime": 0,
hduser@bmse-OptiPlex-3060:~/basic-container-project$
```

Ln 36, Col 29 Spaces:2 UFF:8 LF Compose ⚡ Port:5500 ⚡

## 5. kata

The screenshot shows a Visual Studio Code interface with the following details:

- File Explorer:** Shows files like `index.js`, `docker-compose.yml`, and `Dockerfile`.
- Terminal:** Displays command-line output from a Docker container named `basic-container-project`. The commands run include:
  - `sudo docker compose up -d`
  - `sudo docker inspect mongodb | grep "IPAddress"` (Output: IP address 172.23.0.4)
  - `sudo docker inspect prometheus | grep "IPAddress"` (Output: IP address 172.23.0.4)
  - `sudo docker inspect node-server | grep "Runtime"` (Output: Runtime: kata-runtime)
  - `sudo docker info | grep "Kernel"` (Output: Kernel Version: 4.15.0-189-generic)
- Status Bar:** Shows the current file is `index.js`, line 83, column 15, and the port is 5500.

## Essential Metrics for Container Monitoring

When considering metrics it will be helpful to pick a concrete metrics platform in order to show examples of collection and alerting. In the world of the Cloud Native technologies, Prometheus is the obvious choice for metrics collection, alerting, and in conjunction with Grafana visualization. The four golden signals of monitoring are latency, traffic, errors, and saturation. If one can measure any four metrics of your user-facing system, emphasis should be laid on the four given below.

**Latency** — The time it takes to service a request

**Traffic** — A measure of how much demand is being placed on your system

**Errors** — The rate of requests that fail

**Saturation** — How “full” your service is.

However another method of abstraction regarding docker container metrics is the USE method which we have adopted for our project. “USE” stands for Utilization, Saturation and Errors when looking at the **resources** in a given system.

They have the following definitions:

- **Resource**: all physical server functional components (CPUs, disks, busses, ...)
- **Utilization**: the average time that the resource was busy servicing work
- **Saturation**: the degree to which the resource has extra work which it can't service, often queued
- **Errors**: the count of error events

The most important resources containers provide are CPU, Memory, Network and Disk.

We have queried the relevant metrics in PROMQL. The four metrics given by prometheus are the important ones. Cgroup is a linux feature to limit, police, and account the resource usage for a set of processes. It provides mechanism to limit and monitor system resources like CPU time, system memory, disk bandwidth, network bandwidth, etc. The cgroups works by dividing resources into groups and then assigning tasks to those groups. Docker uses cgroups to limit the system resources.

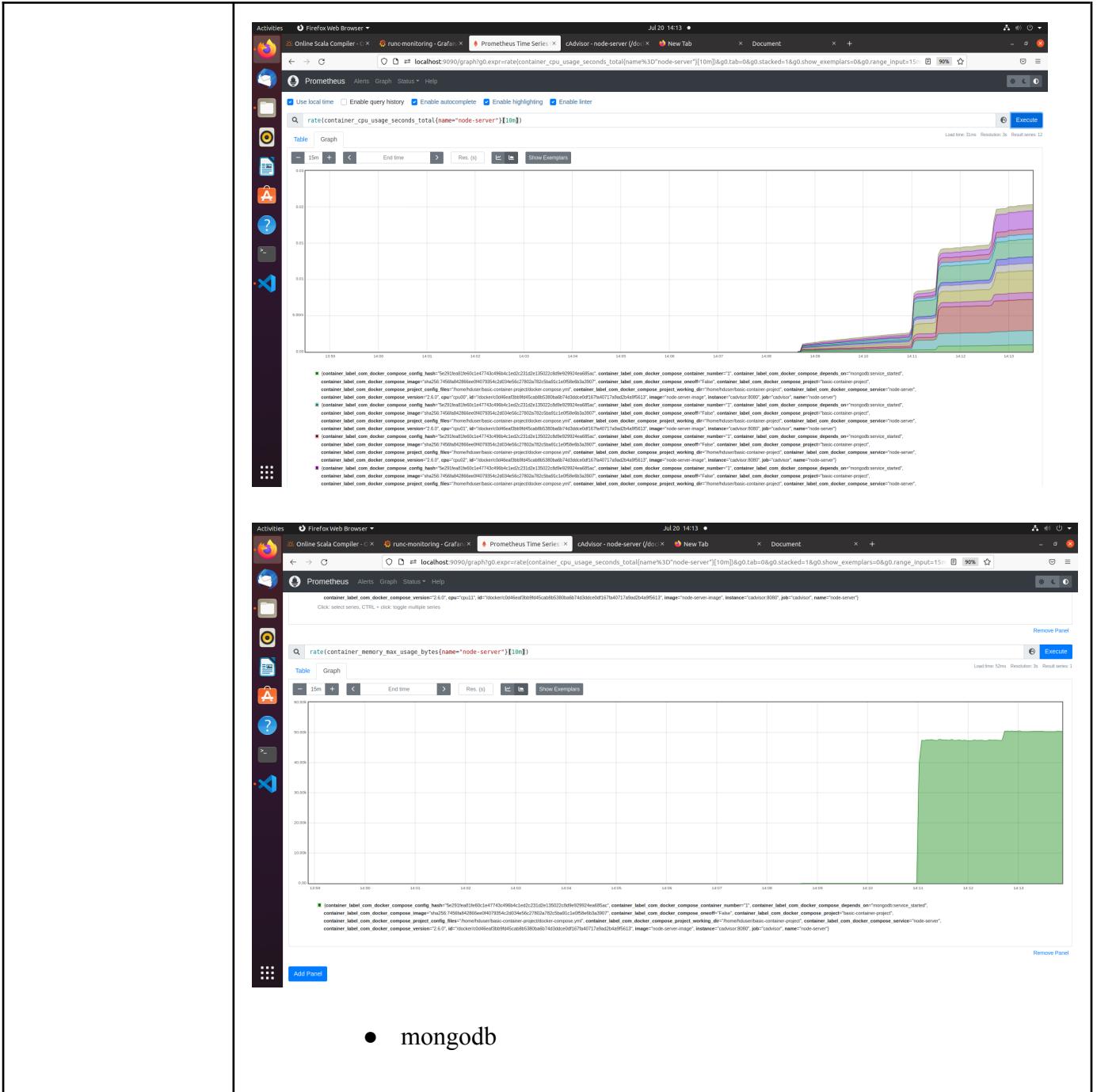
Expression	Description	For
<code>rate(container_cpu_usage_seconds_total{name="container_name"} [1m])</code>	The cgroup's CPU usage in the last minute	The given container
<code>container_memory_usage_bytes{name="container_name"}</code>	The cgroup's total memory usage (in bytes)	The given container
<code>rate(container_network_transmit_bytes_total[1m])</code>	Bytes transmitted over the network by the container per second in the last minute	All containers
<code>rate(container_network_receive_bytes_total[1m])</code>	Bytes received over the network by the container per second in the last minute	All containers

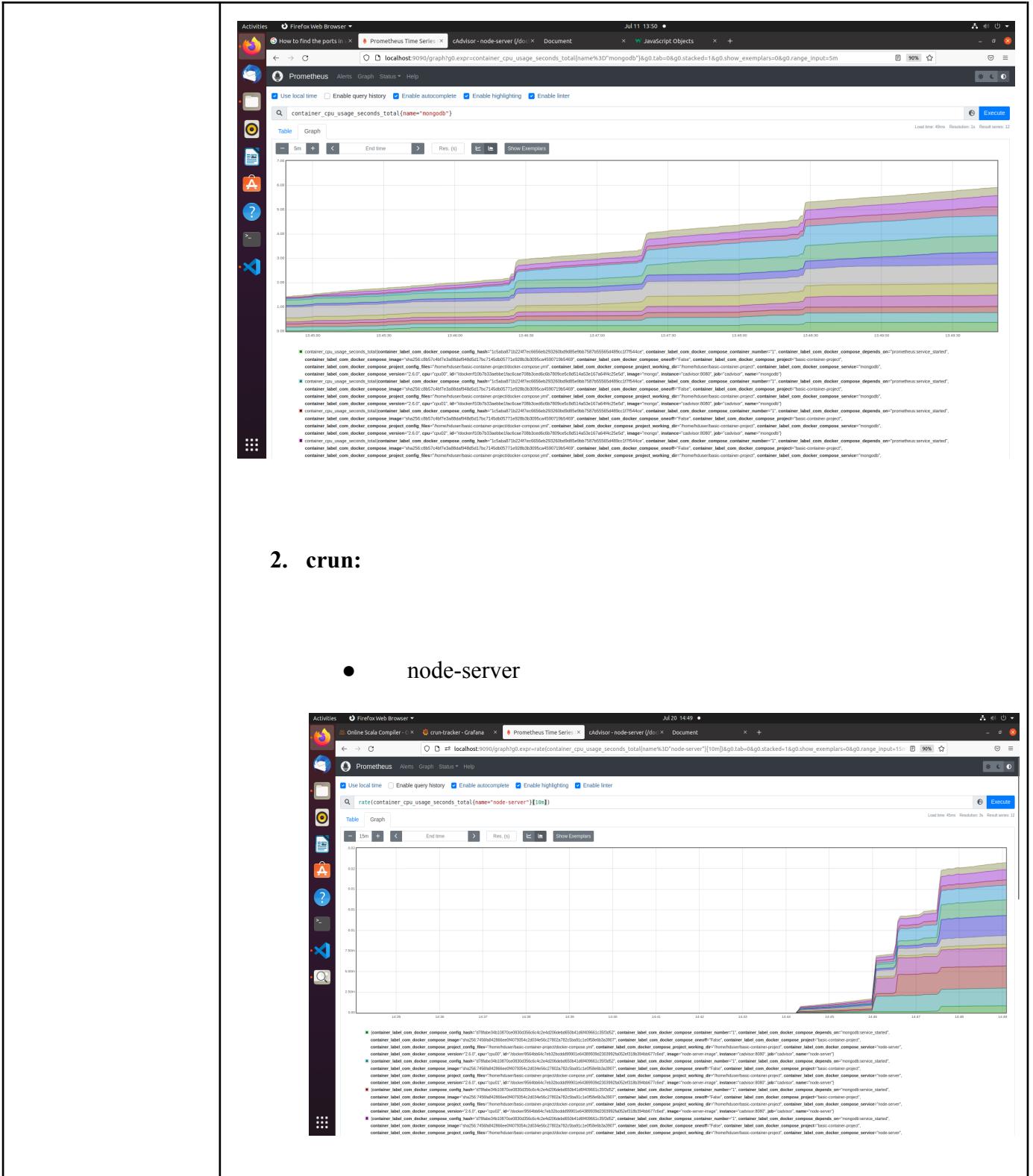
## Prometheus scrapes for one request

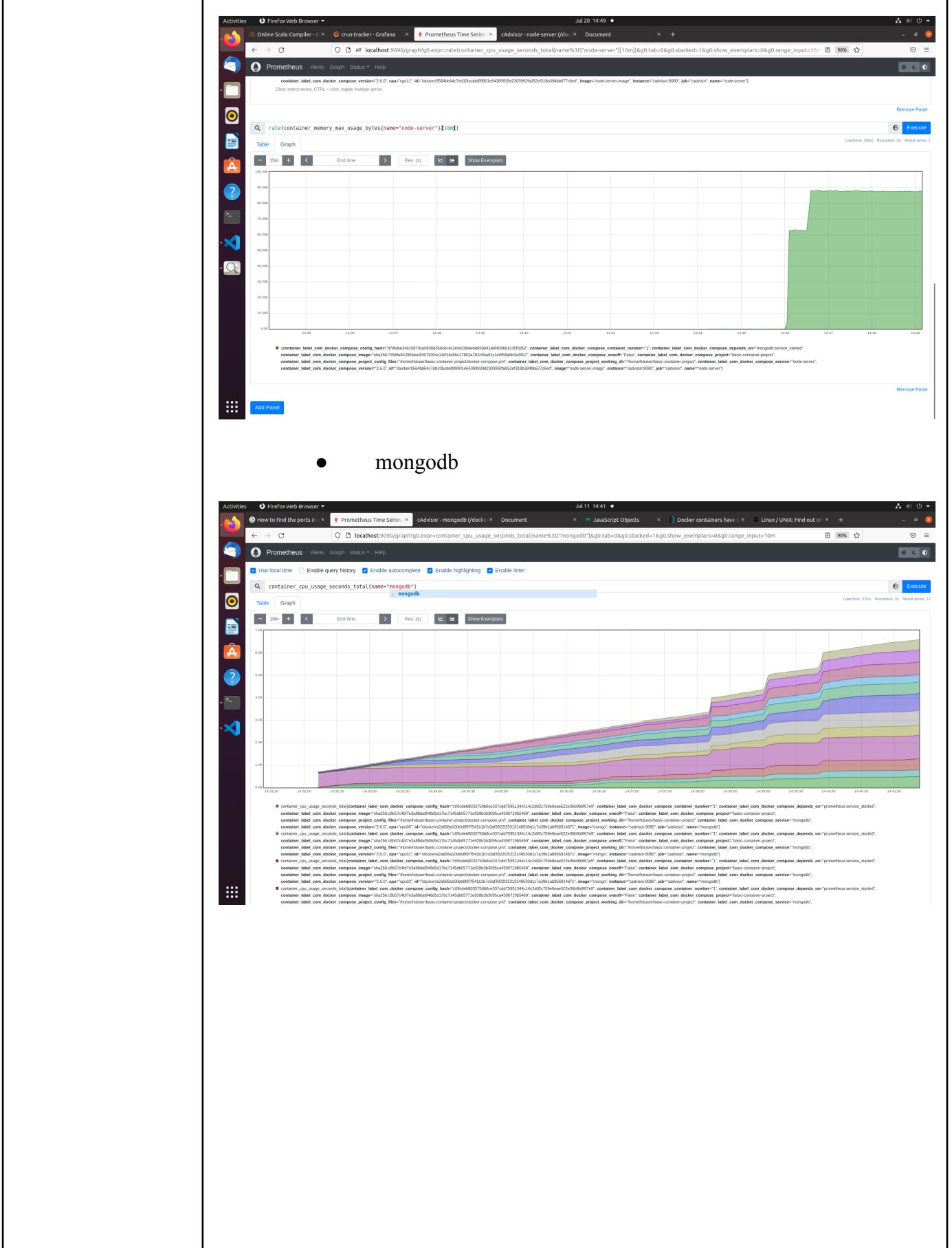
In every case two containers were taken into consideration namely node-server and mongodb

### 1. runc:

- node-server

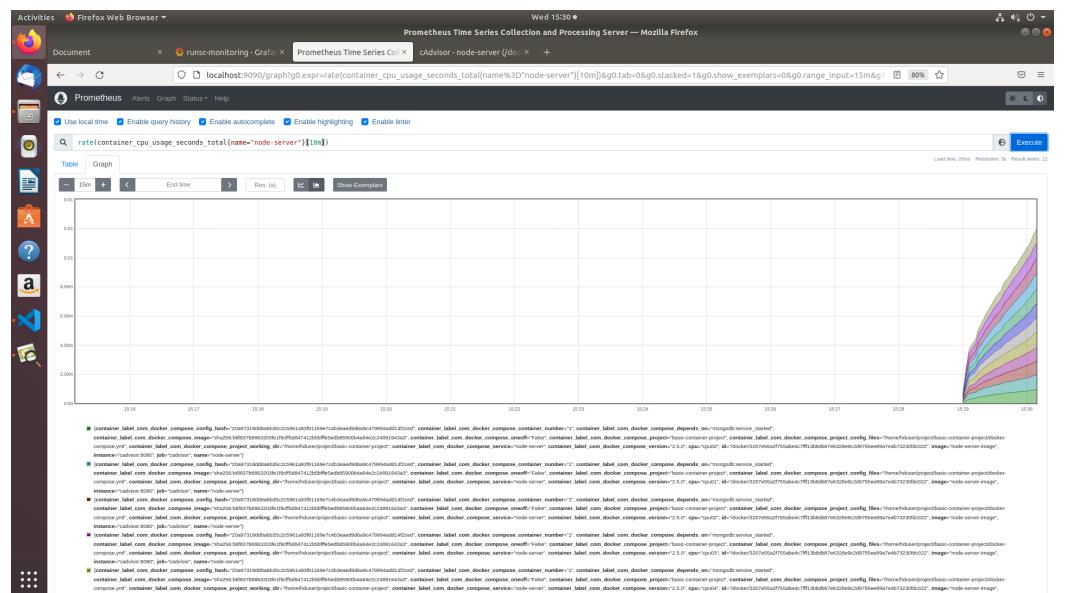
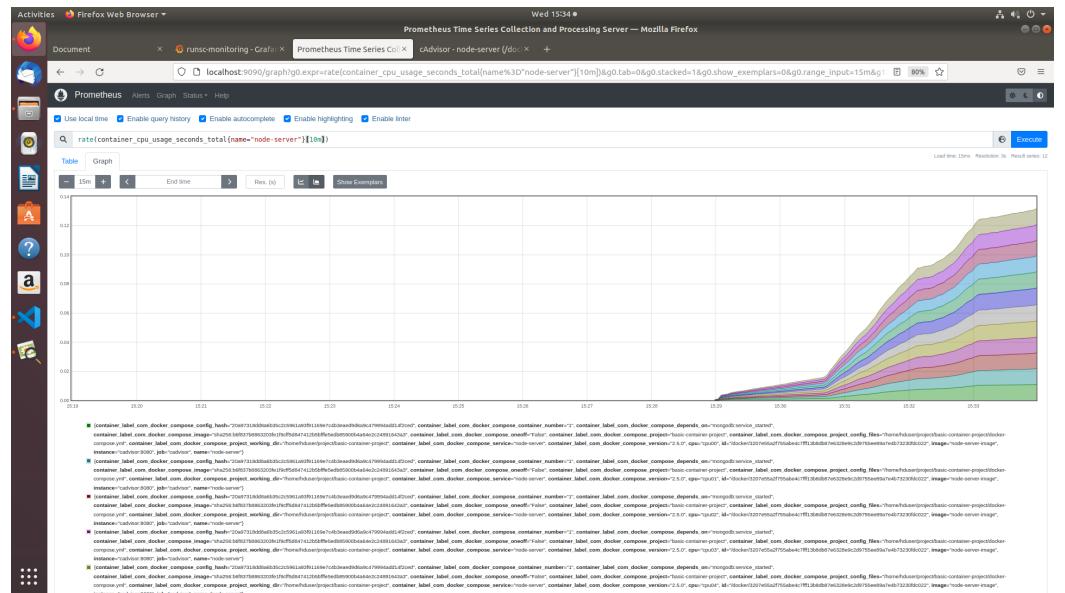






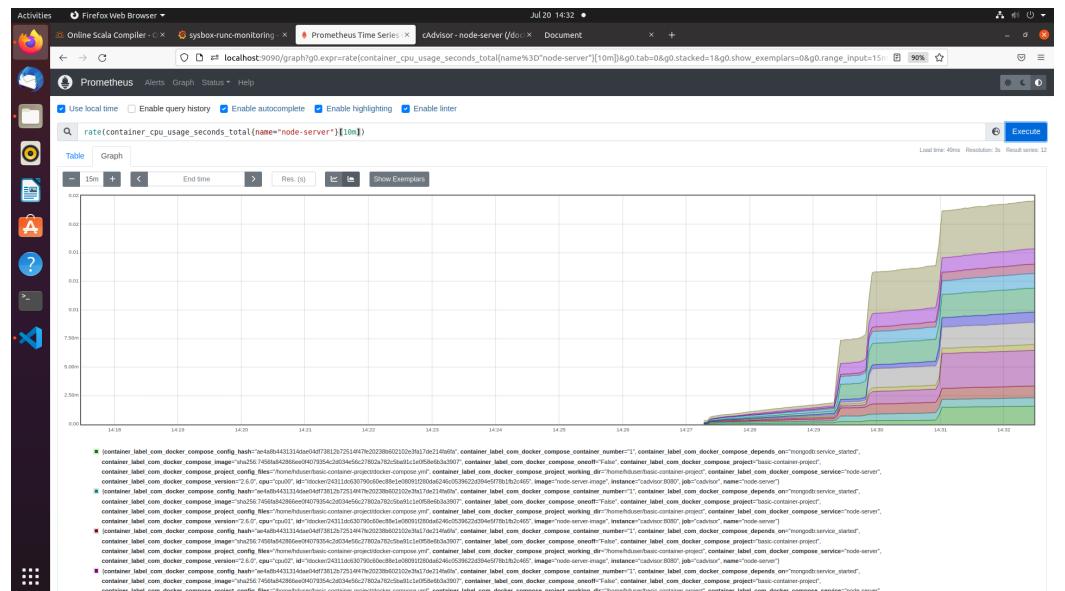
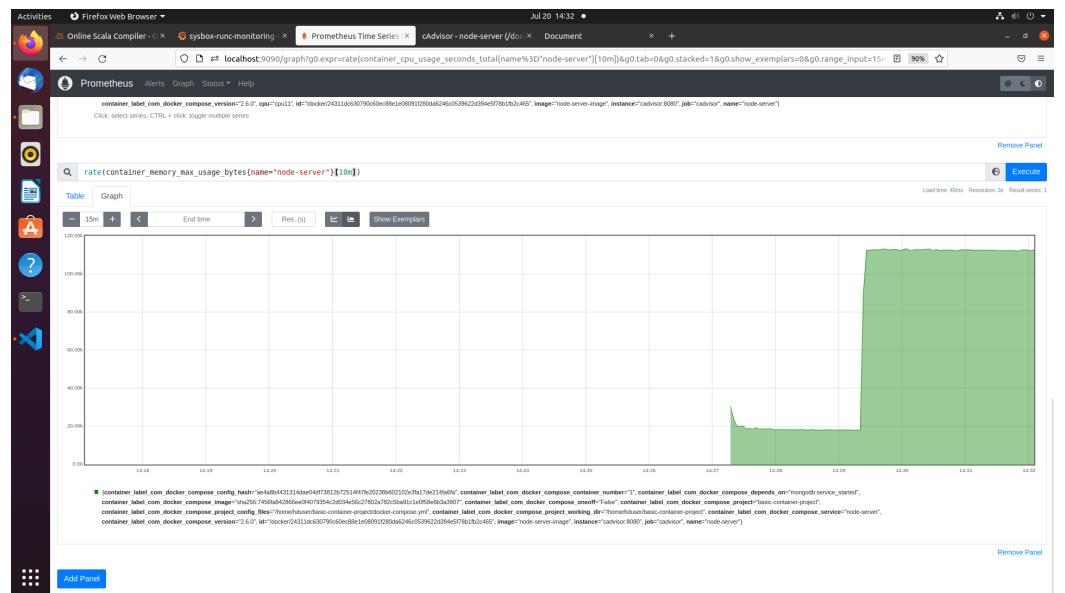
### 3. runsc

- node-server

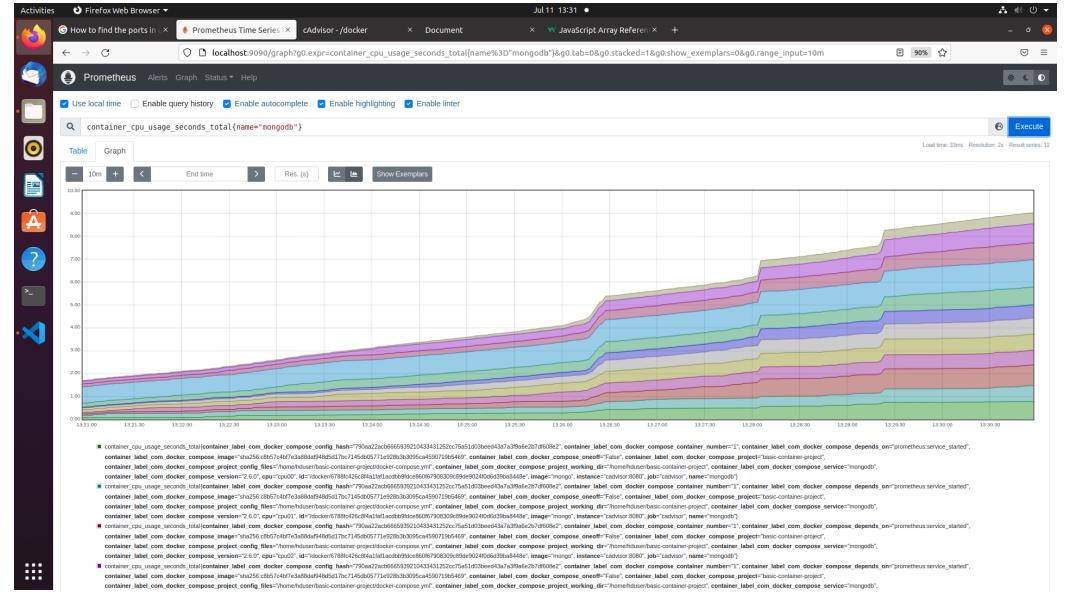


## 4. sysbox-runc:

- node-server

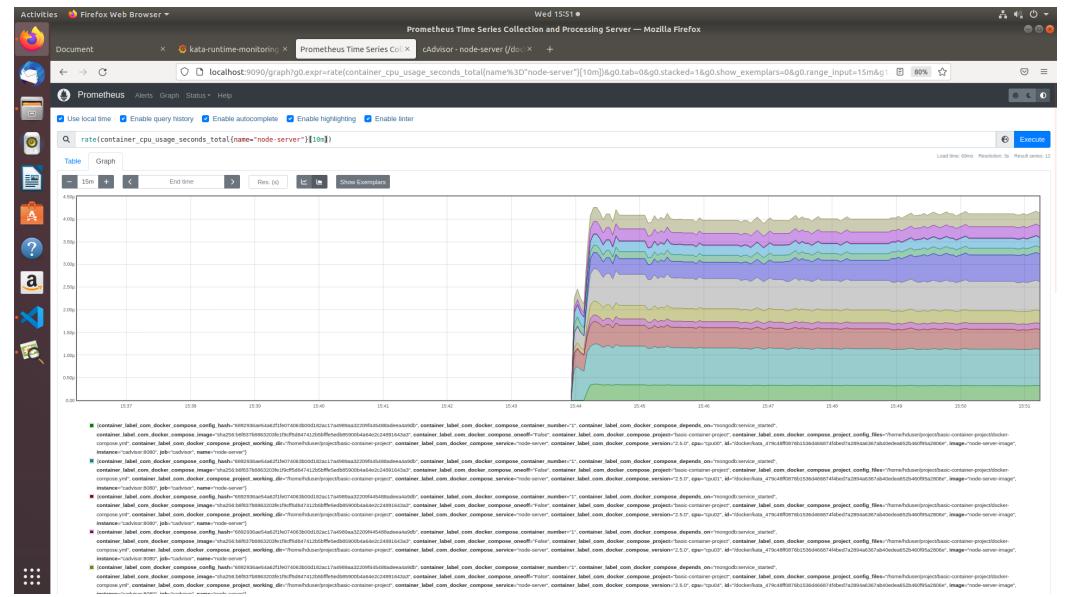


## ● mongodb



## 5. kata:

### ● node-server

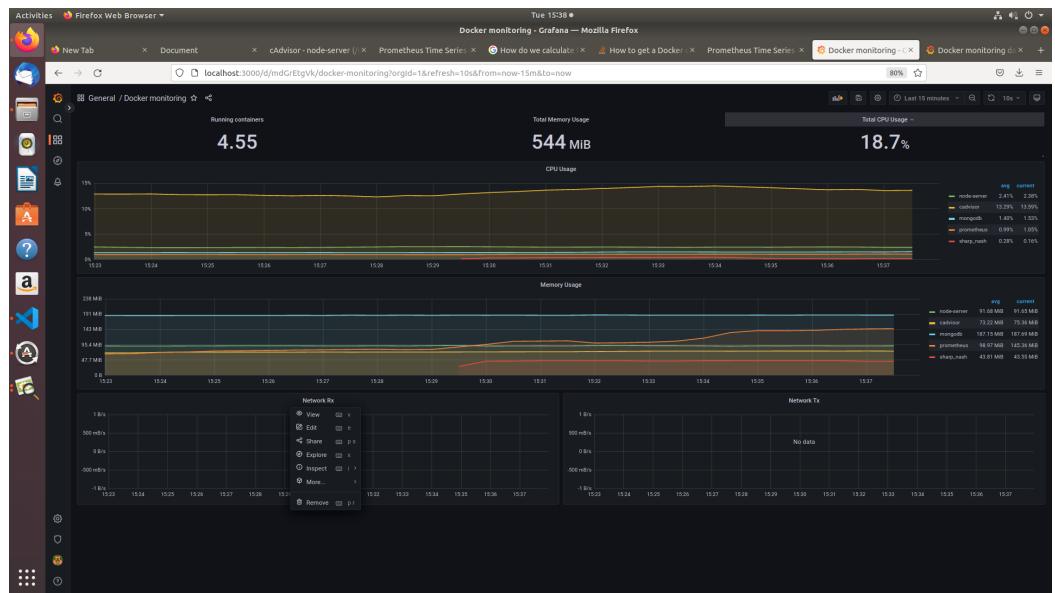


## Grafana dashboard for metrics

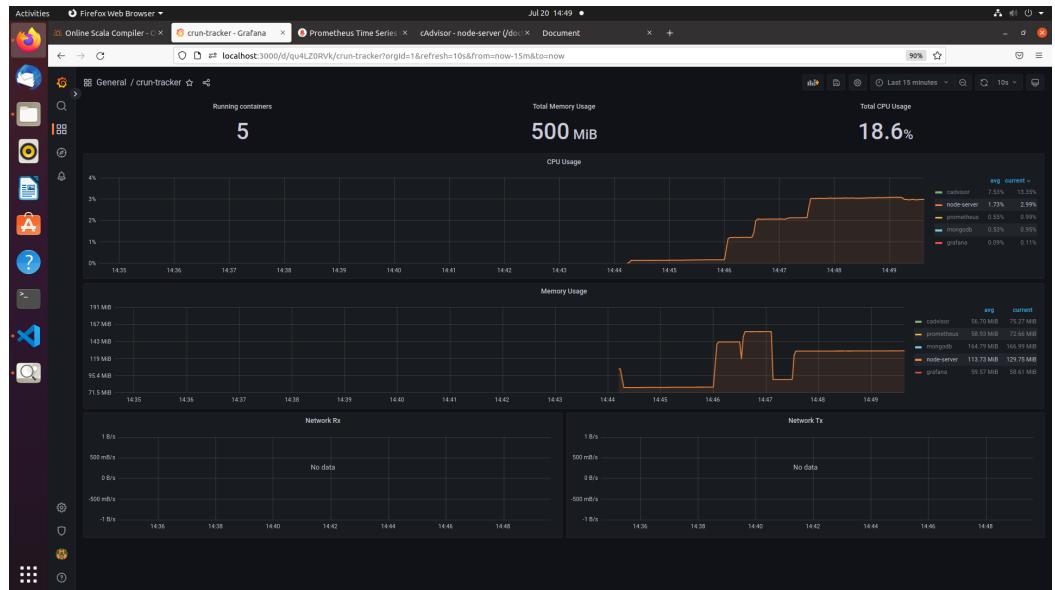
To get a better visualization with cpu usage breakdown and network utilization and to ensure easier comparison we used the standard grafana dashboard for docker monitoring <https://grafana.com/grafana/dashboards/193>. Given below are the screenshots of the dashboards with each runtime.

### Adding a data source:

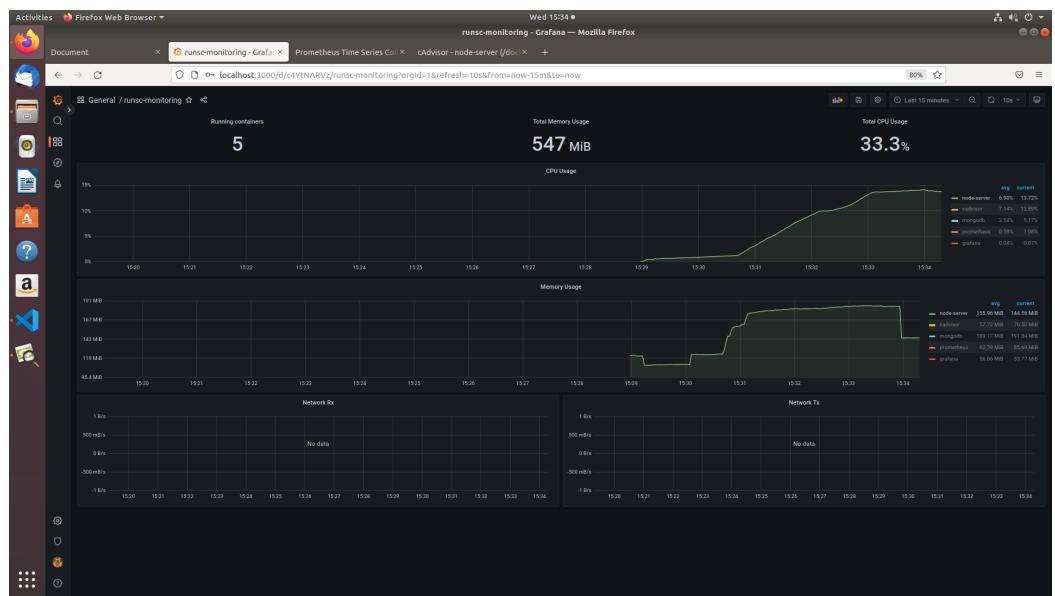
## 1. runc:



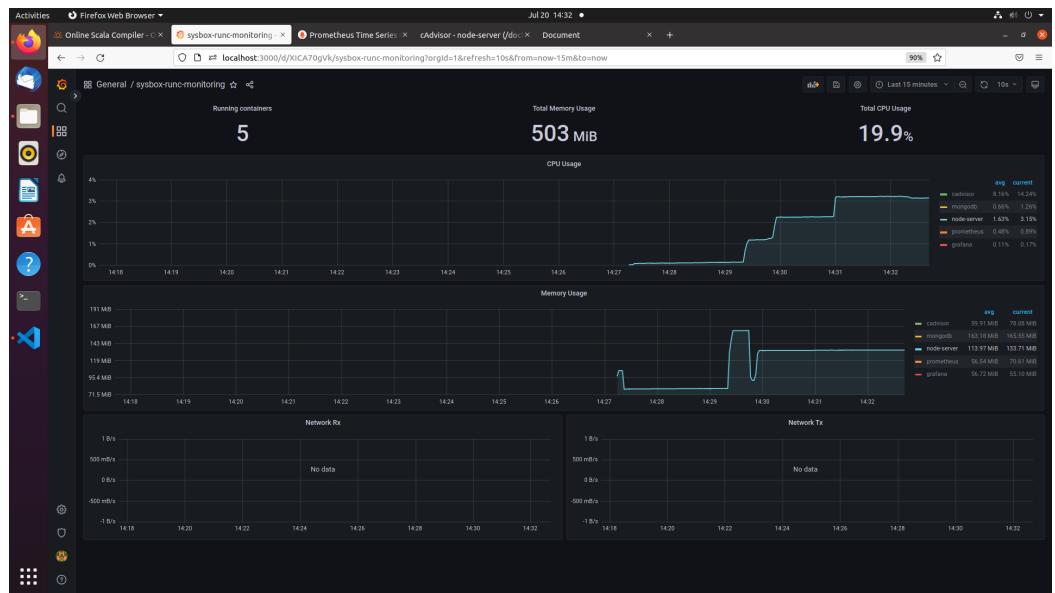
## 2. crun:



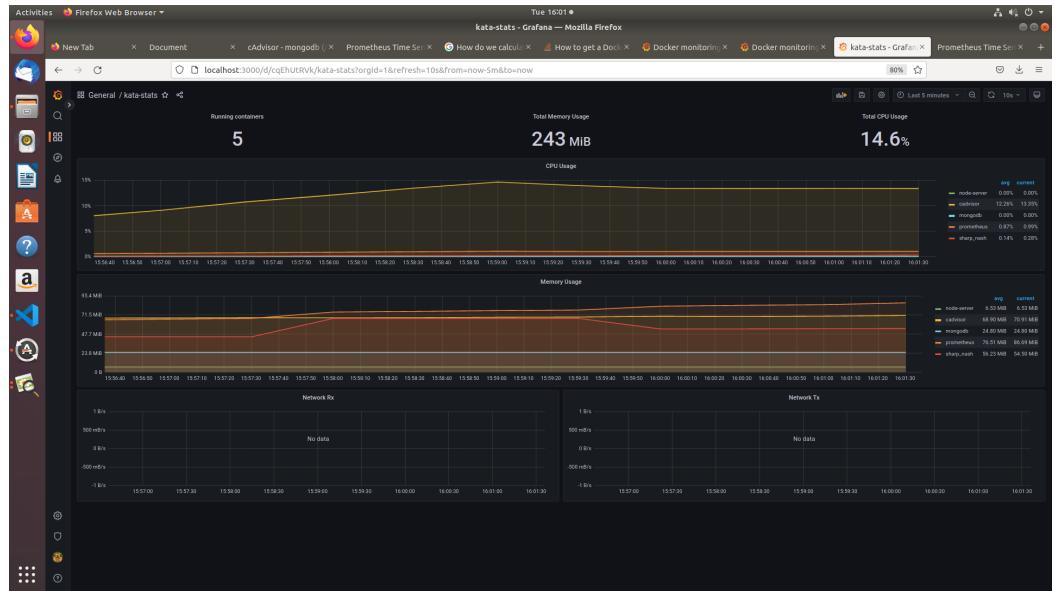
## 3. runsc:



## 4. sysbox-runc:



## 5. kata:



### Interpretation of metrics

The metric used here is “cpu usage”. This is a counter metric that counts the number of seconds the CPU has been running in a particular mode. The CPU has several modes such as iowait, idle, user, and system. Because the objective is to count usage, we use a query that excludes idle time:

```
sum by (cpu)(node_cpu_seconds_total{mode!="idle"})
```

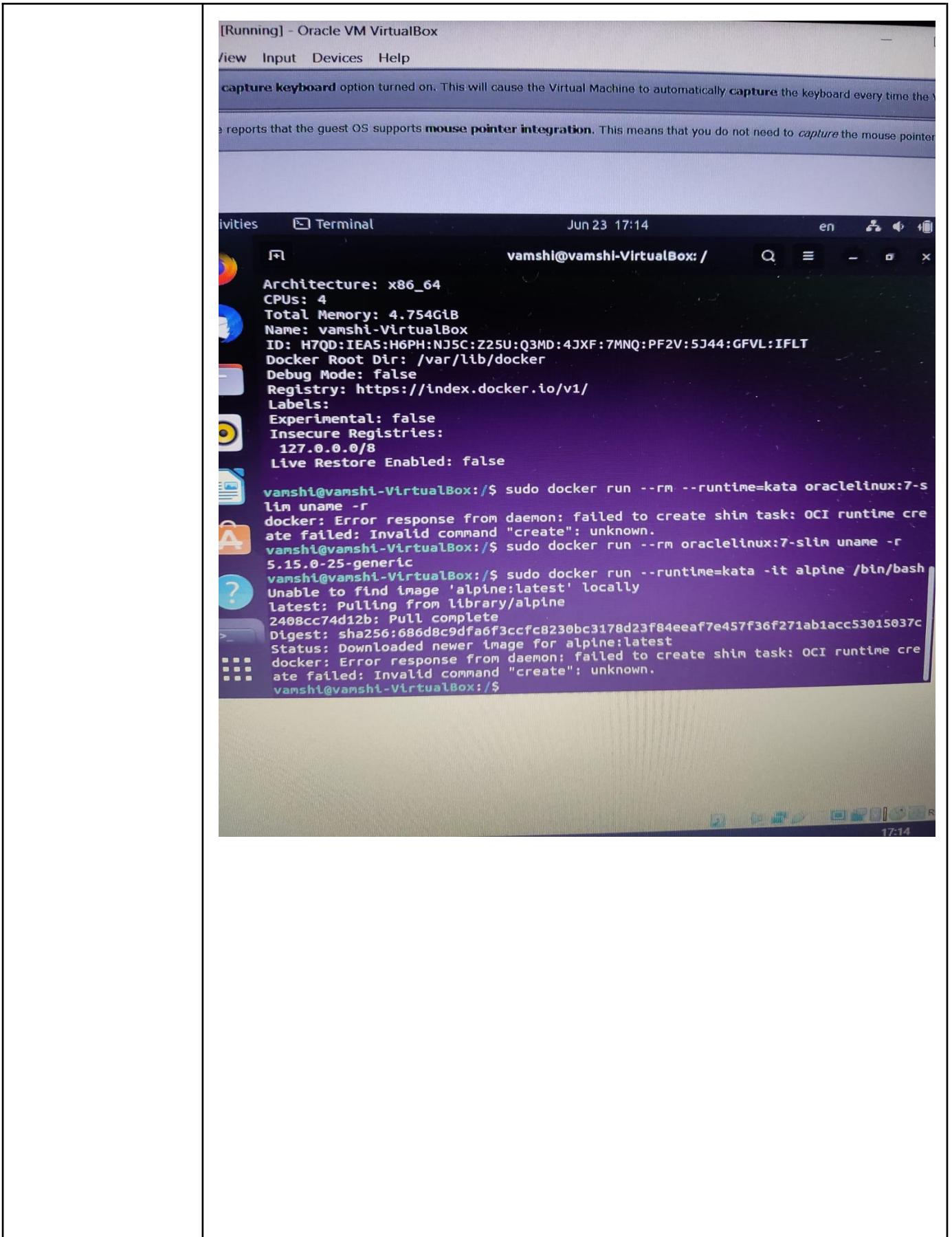
The sum function is used to combine all CPU modes. The result shows how many seconds the CPU has run from the start. To tell if the CPU has been busy or idle recently, use the rate function to calculate the growth rate of the counter:

```
(sum by (cpu)(rate(node_cpu_seconds_total{mode!="idle"})[5m]))*100
```

The above query produces the rate of increase over the last five minutes, which lets you see how much computing power the CPU is using

By default, a container has no resource constraints and can use as much of a given resource as the host's kernel scheduler allows. memory usage by subtracting cache usage from the total memory usage. The API does not perform such a calculation but rather provides the total memory usage and the amount from the cache so that clients can use the data as needed

Out of all the five runtimes kata runtime gives the best cpu and memory utilization without placing any limits on it. Other than kata the default runc runtime gives a good performance along with crun. When looking only at memory utilization for 5 containers over the same time interval of 15 minutes these are the results and Kata could be a suitable open source alternative to runc. However Kata 2.x support for Docker will be removed due to incompatibility of dockershim. as shown below.



<b>Learning Outcomes</b>	<ul style="list-style-type: none"> <li>• Learnt about containerization pertaining to docker in depth</li> <li>• Explored open-source tools like prometheus and grafana for identifying and analyzing relevant metrics for monitoring containers.</li> </ul>
<b>Future Work</b>	<ul style="list-style-type: none"> <li>• Comparing additional OCI runtimes with Kubernetes as most open source runtime's compatibility with docker is depreciating by the end of 2022.</li> <li>• Investigate other metrics made available by Prometheus.</li> </ul>
<b>References</b>	<ol style="list-style-type: none"> <li>1. <a href="https://www.capitalone.com/tech/cloud/container-runtime/">https://www.capitalone.com/tech/cloud/container-runtime/</a></li> <li>2. <a href="https://devopstales.github.io/home/container-runtimes/">https://devopstales.github.io/home/container-runtimes/</a></li> <li>3. <a href="https://github.com/nestybox/sysbox">https://github.com/nestybox/sysbox</a></li> <li>4. <a href="https://github.com/kata-containers/runtime">https://github.com/kata-containers/runtime</a></li> <li>5. <a href="https://github.com/google/gvisor">https://github.com/google/gvisor</a></li> <li>6. <a href="https://blog.freshtracks.io/a-deep-dive-into-kubernetes-metrics-part-3-container-resource-metrics-361c5ee46e66">https://blog.freshtracks.io/a-deep-dive-into-kubernetes-metrics-part-3-container-resource-metrics-361c5ee46e66</a></li> <li>7. <a href="https://opensource.com/article/19/11/introduction-monitoring-prometheus">https://opensource.com/article/19/11/introduction-monitoring-prometheus</a></li> <li>8. <a href="https://tsh.io/blog/grafana-custom-dashboard/">https://tsh.io/blog/grafana-custom-dashboard/</a></li> <li>9. <a href="https://medium.com/aeturnuminc/configure-prometheus-and-grafana-in-dockers-ff2a2b51aa1d">https://medium.com/aeturnuminc/configure-prometheus-and-grafana-in-dockers-ff2a2b51aa1d</a></li> <li>10. <a href="https://docs.docker.com/config/daemon/prometheus/">https://docs.docker.com/config/daemon/prometheus/</a></li> <li>11. <a href="https://jpinjpblog.wordpress.com/2020/04/01/monitoring-and-collecting-docker-container-statistics/">https://jpinjpblog.wordpress.com/2020/04/01/monitoring-and-collecting-docker-container-statistics/</a></li> <li>12. <a href="https://mahesh-mahadevan.medium.com/monitoring-docker-containers-on-windows-using-prometheus-grafana-c32bbb7ed04">https://mahesh-mahadevan.medium.com/monitoring-docker-containers-on-windows-using-prometheus-grafana-c32bbb7ed04</a></li> </ol>