

An Efficient and Scalable Traffic Load Balancing Based on Webserver Container Resource Utilization using Kubernetes Cluster

Ashok L Pomnar¹,

¹Student of ME Computer,
Department of Computer Engineering,
AVCOE Sangamner 422 605,
Maharashtra, India.

Dr. S.K. Sonkar²

Assistant Professor,
Department of Computer Engineering,
AVCOE Sangamner 422 605,
Maharashtra, India

Abstract:- A Today digital transformation era, traditional web based systems and application architecture become the bottleneck of scaling, high availability, portability, and many more downtime and performance-related challenges so the world is started adopting the cloud and cloud-based services to get its benefits. Virtualization technology played a very crucial role to migrate applications from traditional physical systems to virtualized systems but again because of limitations and challenges associated with virtualized system scalability, high availability, and application portability, the Container technique is gaining increasing attention in recent years and has become an alternative to traditional virtual machines to run the application.

Few of the primary motivations for the enterprise to adopt container technology include its convenience to encapsulate and deploy applications, lightweight, as well as efficiency and flexibility in resource sharing. Considering current vertical and horizontal application scaling challenges for the high traffic websites, in this paper we explains the benefits of cloud technology, virtualization, container and Kubernetes clustering features. Also, discuss the algorithm for building dynamic scaling large traffic platforms and applications, which will be runtime, scaled up and down as per user request and traffic demand.

Recent growth in e-business, online web application and mobile user accessibility are exponentially grown with the widespread of Internet of Things, Data Processing and data analytic and many more online portal. To improve better web application performance, web application availability and scalability resource utilization many service provider start design services using micro services and deploy them on the container using Kubernetes cluster. However, the exiting approaches fail to address service availability, handling high traffic load result to dissatisfaction to end users and business impact. To address these issues, in proposed work evaluate many interesting aspects of running high traffic website applications on containerized dynamic scaling and high availability cluster. Proposed Setup run on the cloud again, such as how convenient the execution environment can be set up, what are makes pans of different workloads running in each setup, how efficient the hardware resources, such as CPU and memory, utilized, and how well each environment can scale. The

results show that compared with virtual machines, containers are more easy-to-deploy and scalable environment for high traffic workloads.

Keywords:- Docker, Cloud Computing, Kubernetes Cluster, K8S, Micro service, Web Server, Haproxy Load Balancer, NFS Storage.

I. INTRODUCTION

In today's digital transformation area, Docker and Kubernetes have revolutionized the way of DevOps consulting and both are leading container orchestration tools[1]. There is always a challenge to control an increase in the demand for scaling and auto-healing of the network and virtual instances. Managing the containers is always a task for any company because microservices that are running on the containers do not communicate with each other. They work independently as a separate entity. This is where Kubernetes steps in. Kubernetes is nothing but a platform to manage containers. These containers can be Docker containers or any other alternative containers. Kubernetes orchestrates, manages, and forms a line of communication between these containers [5].

- **The first Section covers Docker:** Containerization Using Docker, Docker for networking, The Docker File, and hosting a web server using Docker.
- **In the Second Section,** we have covered Kubernetes: The Role and Architecture of Kubernetes, basic concepts, features, Kubernetes clusters, scaling and deploying applications, and hosting a Web Server Using Helm.
- **In the third section,** I have planned to deploy the Cloud platform to deploy the complete. Kubernetes and Dockers on Virtual machines which run on either Xen or VMware Cloud. Its compatibility, services, features, and how dockets Kubernetes and Cloud go hand in hand and last.
- **The last section** deploys the back-end script to monitor the traffic on the load balancer server and then automatically scales up and down the webserver containers to distribute the traffic. To summarize, the project will be by taking advantage of each of the cloud and cloud services emerging technology to deploy the highly scalable and high availability platform to run high traffic applications.

II. LITERATURE REVIEW

- Ruchika Muddinagiri, Shubham Ambavane, Simran Bayas [3] In this paper, the author has deployed the containerization application using docker and minikube tools on the local system, so there is future scope to build Kubernetes cloud-deployed scalable application on same which can provide availability as well as scalability.
- RobertBotez,Calin-MarianIurian,Iustin-AlexandruIvanciu,VirgilDobrota [5]: In this paper, authors have evaluated a solution for monitoring vehicles in real-time using containerized applications. However, there is the future scope for work on evaluating the Docker performance and security as well set up a Kubernetes cluster on a private cloud. The cluster nodes are made up of multiple Raspberry Pi platforms that will collect more sensor information.
- Jay Shah, Dushyant Dubaria [4]: In this paper, the WordPress blog application is deployed using containerization technology with a standalone and static resource on the Google Cloud Platform. The major advantage of using Kubernetes is orchestrating between many applications. It is also very useful for scaling many applications in very less time. By using this, we can save our costs and time.

III. PROPOSED METHODOLOGY

A. Architecture

In this paper, three virtual Machines need to deploy the Virtual Machin on the Cloud Service Provider (ESDS eNlight360/AWS/Microsoft Azure) with high available Cloud. Install and set up the Kubernetes cluster on all these three-node with one master and 2 worker/slave kinds of architecture where we will get the benefit of container failover advanced in case any VM is failed/rebooted.

We can deploy any PHP application container using apache-PHPDocker container using Kubernetes deployment,pods, services, and endpoints features and make the application accessible with IP address and port number using Kubernetes networking. A persistent volume is used to save and store application data after failover of VM, Kubernetes nodes, or container restarting.

We can deploy another HaProxy container to distribute the traffic of web server traffic among different application containers, which are dynamically scaling up and down as container resource utilization and traffic increase. Kubernetes scheduler and monitoring continuously monitor the traffic utilization and CPU and memory utilization of the container and increasing application container using application deployment, which quickly gets available behind the HaProxy load balancer to server the web traffic.

When web request, CPU, and memory utilization are decreased, the Kubernetes scheduler and Kubernetes manager are reducing the container and back to the default min container count.

We do use the Round Robin algorithm to equally balance the load among all application containers.

B. Main System Components

- **Centos VM:** Setup the three CentOS-based Virtual machines on any Cloud platform.
- **Setup and Configuration Kubernet Cluster:** Setup and configure the three-node one Master node and 2 Worker-Slave nodes Kubernetes cluster on these three VM.
- **Deploy Application Container:** Deploy the apache-PHP web servers, MySQL database Docker container, and installed on the Kubernetes cloud.
- **Application Accessibility:** With help of Kubernetes networking, services and endpoint make the application accessible with a web serverURL.
- **Web Server Load Balancer:** Deploy the haproxy Docker container on Kubernetes cloud with dynamical container scaling feature automatically detect to enable or disable the application container whenever traffic demand is increasing or decreasing.
- **Web Traffic generating tool:** Using the Siege or curl method to generate real-time traffic on the webserver.

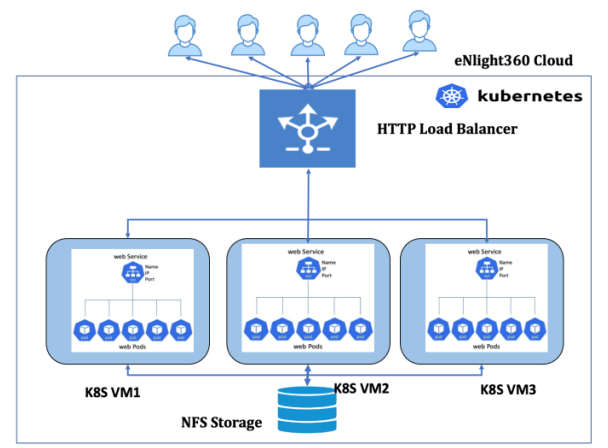


Fig. 1: Solution Architecture Diagram

The proposed system will be comprised of the following modules:

- **Module 1: Setup phase:** Build and set up the Virtual machine on the cloud to set up the Kubernetes clustering with required resources.
- **Module 2: Build Kubernetes Cloud and Containers:** Build and set up the Kubernetes clustering with high availability option.
- **Module 3: Deploy Application:** Deploy containerized services like web container, database container, and load balancer container to run the WordPresscontainerized application.
- **Module 4: Deploy backend scripts:** Deploy the proposed algorithm in the form of script, which monitor real-time traffic patterns, and as per container traffic handling limit/threshold set it will dynamically scale up, scale down the web servers container, and add the same behind load balance too to balance the traffic.
- **Module 5: Dynamic scaling and High availability:** Demonstrate the container scale up whenever the traffic spike or increased and scale down when traffic is reduced. Also, demonstrate the application’s high availability in terms of any web server container failure, any

Kubernetesnode (VM) failure, VM failure, or any cloud node failure.

- **Module 6: Analysis Reports:** Prepare the details analysis report to consider the traffic pattern, traffic balancing among several web servers, container scale up and scale down as per traffic pattern, and high availability report in terms of components failure.

IV. ALGORITHMS

A. Web Request Balancing Algorithm: Round robin

We do use the Round Robin algorithm to balance the load among all application containers. The round-robinload-balancing algorithm is one of the popular methods for distributing web traffic across a group of web servers. Going down the list of servers in the group, the round-robin load balancer forwards a request to each server in turn. When it reaches the end of the list, the load balancer loops back and goes down the list again (sends the next request to the first listed server, the one after that to the second server, and so on).

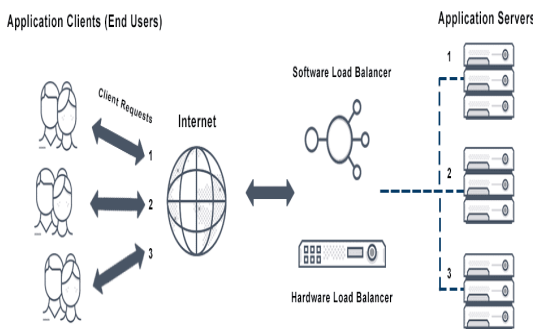


Fig. 2: Web Request balancing using Round Robin Algorithm

B. Horizontal Container Scaling Algorithm:

Desired App Replicas = $\text{ceil}[\text{current App Replicas} * (\text{current Scaling Metric Value} / \text{desired Scaling Metric Value})]$

Desired App Replicas: Application Replicas pod count that will sent to the controller after calculations.

Ceil (): This function that rounds a fractional number. For example, $\text{ceil}(11.28)$ is 12.

Current App Replicas: Current number of Application Replicas pods for a given deployment or any other superset of “scale” object type.

Current Scale Metric Value: Current value of metric for a given scaling factor metric. Can be 800m or 1.5Gi, for custom metrics, it can be 500 events per second, etc.

Desired Scale Metric Value: Metric that has been set to maximum scalability of application pod. Eventually, with all mechanisms provides, your app runs at this metric value. This value should not be too low or too high.

Let us say we have a scaling configuration with a target CPU usage of 70%, a minimum pod count of 4, and a maximum pod count of 20.

The current deployment status is six pods averaging 95% usage.

$$\text{Desired Replicas} = \text{ceil}[6 * (85/70)] = \text{ceil}(8.14) = 9$$

Scaling down Algorithm:

Let us say we have a pod-scaling configuration with a target CPU usage of 70%, a minimum pod count of 4, and a maximum pod count of 20.

The current application pod deployment status is: There are 10 total pods. 10 pods averaging 45% usages.

The same algorithm we used to scale down the application pod:

The formula is applied to all normal pods.

$$\text{Desired Replicas} = \text{ceil}[10 * (45/70)] = \text{ceil}(6.42) = 7 > 10$$

With the above calculation three application pods are extra, so as per the scaling downtime set in the configuration it will scale down the application pod to seven.

V. WORKING

A. Kubernetes Cloud and Container Status:

We can set up and create three Centos VM virtual Machines on the eNlight360 Cloud to build the High Available Kubernetes Master-Slave clustering. Below is the Kubernetes cluster node status

```
[root@k8snode1 ~]# kubectl get node
NAME                                STATUS    ROLES    AGE    VERSION
k8snode1.alexparkar.com            Ready    control-plane,master    201d    v1.22.2
k8snode2.alexparkar.com            Ready    <none>    201d    v1.22.2
k8snode3.alexparkar.com            Ready    <none>    195d    v1.22.2
[root@k8snode1 ~]#
```

Fig. 3: Kubernetes cluster status

B. Deploy the apache-PHP base application container:

We can deploy and install the PHP-Apache web server application container using k8s deployment components which run on the respective worker nodes with a minimum number of the replica set and maximum, minimum replica scaling number when pod CPU and memory demands increased because of a total number of web request increases.

```
[root@k8snode1 ~]# kubectl get pod,deploy -n mydemo
NAME                                READY    STATUS    RESTARTS    AGE
pod/centos-loadtest-966f54885-8g4mx  1/1     Running   1 (161d ago)  194d
pod/mydemo-app-6bf8bdc664-f29hb      1/1     Running   0            28d

NAME                                READY    UP-TO-DATE    AVAILABLE    AGE
deployment.apps/centos-loadtest      1/1     1             1            194d
deployment.apps/mydemo-app           1/1     1             1            194d
[root@k8snode1 ~]#
```

Fig. 4: K8s pod, deploy status

VI. CONCLUSION

By adopting the cloud native services benefit, Virtualization, containerization technology are used to build and migrate applications from a traditional physical system to virtualized system, the containerized system which convenient to encapsulate and deploy applications, lightweight operations, as well as efficient and flexible in resources scaling.

Considering current vertical and horizontal application scaling challenges for high traffic websites, in this paper taking the benefit of the cloud, virtualization, container, and Kubernetes clustering features we proposed a dynamic scaling algorithm considering backend latency for large traffic platforms and applications which will be dynamically scaled up and down as per user request and traffic demand.

REFERENCES

- [1.] Ruchika Muddinagiri, Shubham Ambavane, Simran Bayas , "Self-Hosted Kubernetes: Deploying Docker Containers Locally With Minikube" , IEEE Xplore, Conference,18 August 2020 ISBN Information: ,DOI: 10.1109/ICI- TAET47105.2019.9170208"
- [2.] Jay Shah,Dushyant Dubaria",Building Modern Clouds: Using Docker, Kubernetes Google Cloud Platform",IEEE Xplore Conference: 14 March 2019 ,DOI: 10.1109/CCWC.2019.8666479
- [3.] Robert Botez,Calin-Marian Iurian,Iustin-Alexandru Ivanciu,Virgil Dobrota ,,"Deploying a Dockerized Application With Kubernetes on Google Cloud Platform", IEEE Xplore Conference: 16 July 2020,DOI: 10.1109/COMM48946.2020.9142014
- [4.] K. Matthias, and S.P. Kane, "Docker:Up and Running", O'Reilly, 2015
- [5.] Deploy on Kubernetes Documentation , Aug 2021 [Online], Available: <https://docs.docker.com/desktop/kubernetes/>
- [6.] Kubernetes Cloud Setup Documentation, Aug 2021, [Online] Available: <https://kubernetes.io/docs/setup/production-environment/>
- [7.] B. Burns, B. Grant, D. Oppenheimer, E. Brewer, J. Wilkes, "Borg, omega, and kubernetes", 2016.
- [8.] Q. Zhang, L. Liu, C. Pu, Q. Dou, L. Wu and W. Zhou, "A Comparative Study of Containers and Virtual Machines in Big Data Environment," 2018 IEEE 11th International Conference on Cloud Computing (CLOUD), San Francisco, CA, 2018, pp. 178-185, doi: 10.1109/CLOUD.2018.00030.
- [9.] 'DevOps; Puppet, Docker and Kubernetes – Learning path' by Thomas Uphill, Arundel, Khare, Saito, Lee and Carol Hsu, Packt Publications, First Edition, 2017
- [10.] 'Cloud Native Applications- The Intersection of Agile Development and Cloud Platforms' by Douglas Bourgeois, David Kelly, Thomas Henry members of Deloitte Touche Tohmatsu Limited,2016
- [11.] 'Kubernetes from the ground up, deploy and scale performant and reliable containerized applications with Kubernetes' by Level Up Kubernetes Program, Basit Mustafa, Tao W, James Lee, Stefan Thorpe, 2018