



e-ISSN: 2319-8753 | p-ISSN: 2347-6710

IJIRSET

International Journal of Innovative Research in
SCIENCE | ENGINEERING | TECHNOLOGY



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

Volume 10, Special Issue 2, December 2021

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA

Impact Factor: 7.569

Cloud Application Infrastructure Autoscaling

Mrs.N.S.Abinaya¹, Dr.M.Deepamalar², C.Laavanya³

Assistant Professor, Paravathy's Arts and Science College, Dindigul, TamilNadu, India¹

Associate Professor, Paravathy's Arts and Science College, Dindigul, TamilNadu, India²

PG Scholar, Paravathy's Arts and Science College, Dindigul, TamilNadu, India³

ABSTRACT: In Cloud Computing, scaling is the process of adding or removing compute, storage, and network services to meet the demands a workload makes for resources in order to maintain availability and performance as utilization increase. Autoscaling also referred to as autoscaling, autoscaling and sometimes automatic scaling is a cloud computing technique for dynamically allocating computational resources. Depending on the load to a server farm or pool, the number of servers that are active will typically vary automatically as user needs fluctuate. Autoscaling and load balancing are related because an application typically scales based on load balancing serving capacity. In other words, the serving capacity of the load balancer is one of several metrics (including cloud monitoring metrics and CPU utilization) that shapes the autoscaling policy.

KEYWORDS:-Cloud Computing- Autoscaling, Load Balancing ,CPU utilization.

I. INTRODUCTION

In cloud computing, Autoscaling offers the facility to the clients to scale up and scale down the resources as per the clients demands. As everything is taking place in automatic manner, so human intervention errors are less and reduce the manpower and cost. The CPU utilization defines the workload of client, which are being transferred across the cloud platform through a proper network with maximum bandwidth. The Google Cloud Platform suite of services is always evolving. Google frequently & periodically introduces changes or may discontinue the service based on user demand or competitive pressures. Google's main competitors in the public cloud computing market include Amazon Web Services (AWS) and Microsoft Azure. What you might be used to thinking of as software and hardware products become services. These services provide access to the underlying resources. The list of available Google Cloud Services is long and it keeps growing. When you develop your website or application on Google Cloud you mix and match these services into combinations that provide the infrastructure you need and then add your code to enable the scenarios you want to build. Google Cloud consists of a set of physical assets, such as computers and hard disk drives and virtual resources such as virtual machines (VMs) that are contained in Google's data centers around the globe. Some resources can be accessed by other resources, across regions and zones. These global resources include preconfigured disk images, disk snapshots and networks. Some resources can ne accessed only be resources include static external IP addresses. Other resources can be accessed only by resources that are located in the same zone. These zonal resources include VM instances their types and disks.

II. METHODS

GOOGLE CLOUD VIRTUAL MACHINE

A virtual machine (or "VM") is an emulated computer system created using software. It uses physical system resources, such as the CPU, RAM, and disk storage, but is isolated from other software on the computer. It can easily be created, modified, or destroyed without affecting the host computer. Virtual machines provide similar functionality to physical machines, but they do not run directly on the hardware. Instead, a software layer exists between the hardware and the virtualmachine. These applications allow you to run multiple VMs on a single computer. For example, Parallels Desktop for Mac allows you to run Windows, Linux, and macOS virtual machines on your Mac. Virtual machines run on a physical machine and access computing resources from software called a hypervisor. The hypervisor abstracts the physical machine's resources into a pool that can be provisioned and distributed as needed, enabling multiple VMs to run on a single physical machine.

INSTANCETEMPLATE

Instance Templates define the look of every virtual machine in the cluster (disk, CPUs, memory, etc). Managed Instance Groups instantiate a number of virtual machine instances using the Instance Template. An instance template is a resource that you can use to create VM instances and managed instance groups. Instance templates define the machine type, boot disk image or container image, labels, and other instance properties. You can then use an instance template to create a managed instance group or to create individual VM instances.

Instance templates are a convenient way to save a VM instance's configuration so you can use it later to create new VM instances or groups of VM instances. An instance template is a resource that you can use to create virtual machine (VM) instances and managed instance groups (MIGs). Instance templates define the machine type, boot disk image or container image, labels, and other instance properties. You can then use an instance template to create a MIG or to create individual VMs. Instance templates are a convenient way to save a VM instance's configuration so you can use it later to create VMs or groups of VMs.

INSTANCEGROUP

An instance group is a collection of virtual machine (VM) instances that you can manage as a single entity. A managed instance group (MIG) contains identical instances that are based on an instance template. Managed instance groups maintain high availability of your apps by proactively keeping your instances available, that is, in the running state. Managed instance groups support autohealing, load balancing, autoscaling, and auto-updating.

Auto scaling is a feature of managed instance groups (MIGs). A managed instance group is a collection of virtual machine (VM) instances that are created from a common instance template. An autoscaler adds or deletes instances from a managed instance group based on the group's autoscaling policy. Although Compute Engine has both managed and unmanaged instance groups, only managed instance groups can be used with an autoscaler.

AUTOSCALING

Google Cloud platform offers autoscaling capabilities that allow us to automatically add or delete instances from a managed instance group based on increase or decrease in the load. Google Cloud Autoscaling helps our applications to gracefully handle the increased traffic and reduces the cost when the need for resources is lower. We just elucidate the autoscaling policy and the autoscaler performs the automatic scaling, based on the measured load. Autoscaling group does not scale out for alarms or scheduled actions that occur while the process is suspended.

Managed instance groups (MIGs) offer autoscaling capabilities that let you automatically add or delete virtual machine (VM) instances from a MIG based on increases or decreases in load. Autoscaling helps your apps gracefully handle increases in traffic and reduce costs when the need for resources is lower. You define the autoscaling policy and the autoscaler performs automatic scaling based on the measured load and the options you configure.

Autoscaling works by adding more VMs to your MIG when there is more load (scaling out, sometimes referred to as scaling up), and deleting VMs when the need for VMs is lowered (scaling in or down).

CLOUD SCALING

We can autoscale based on the average CPU utilization of a managed instance group. The autoscaler collects the CPU utilization of the instances in the group and determine whether it needs to scale. We set the target CPU utilization the autoscaler should maintain and the autoscaler will work to maintain that level.

The autoscaler calculates the target CPU utilization level as a fraction of the average use of all vCPUs over time in the instance group. If the average usage of your total vCPUs is more than the target utilization, the autoscaler will add more virtual machines. For example, If we set target utilization as 0.75, autoscaler tries to maintain an average usage of 75% among all vCPUs in the instance group.

The autoscaler continuously collects usage information based on the policy, compares actual utilization to your desired target utilization, and determines if the group needs to be scaled up or down. For example, if you scale based on CPU Utilization, you can set your target utilization level at 80% and the autoscaler will poll for the CPU utilization and check if it is > 80%, if so it adds the new instance to the instance group, if CPU utilization < 80 then it removes an instance from the group, making sure the capacity is maintained. The autoscaler collects information based on the policy. Then it will compare it to the desired target utilization, and determine if it needs to perform scaling. Scaling up and scaling down refer to two dimensions across which resources and therefore capacity can be added.

Scale Up (Vertical Scaling):

Scaling up is the process of resizing a server (or replacing it with another server) to give it supplemental or fewer CPUs, memory, or network capacity.

Benefits of Scaling Up:

1. Vertical scaling minimizes operational overhead because there is only one server to manage. There is no need to distributed the workload and coordinate among multiple servers.
2. Vertical scaling is best used for applications that are difficult to distribute. For example, when a relational database is distributed, the system must accommodate transaction that can change data across multiple servers. Major relational databases can be configured to run on multiple servers, but its often easier to verticallyscale.

Scale Down (Horizontal Scaling):

Instead of resizing an application to a bigger server, scaling down splits the workload across multiple servers that work in parallel.

Benefits of Scaling Down:

1. Applications that can sit within a single-machine like many websites are well suited to horizontal scaling because there is little need to coordinate tasks between servers. For example a retail website might have peak periods, such as around the end of year holidays. During those time, additional servers can be easily committed to handle the additionaltraffic.
2. Many front-end applications and microservices can leverage horizontal scaling. Horizontal scaled applications can adjust the number of servers in use according to the workload demandpatterns.

LOADBALANCING

A load balancer improves the distribution of work load by redirecting the request to multiple nodes. Load balancing allows to maximize throughput, minimize response time, and increase reliability and availability of the resources. Let us configure a HTTP Load balancer in Google Cloud Platform (GCP) for a single Region .Load Balancers, and especially Cloud Load Balancers can distribute your users among multipleservers. If your visitors grow, the Balancer will auto- scale your resources. At the same time, a Balancer could direct users to specific servers based on their geographical location. Another feature is the capacity of load balancing traffic through an HTTPS connection.

On the basis of functionality, load balancers are classified as hardware load balancer and elastic load balancer. Hardware load balancers are concerned with the distribution of workload at hardware level i.e. memory, storage and CPU.

Elastic Load Balancing automatically distributes incoming application traffic across multiple targets, such as Amazon EC2 instances, containers, and IP addresses. It can handle varying load of user application traffic in a single availability zone or across multiple availability zones.

Elastic Load Balancing offers three types of load balancers that all feature high availability, automatic scaling, an robust security necessary to make user applications fault tolerant.

Application Load Balancer operates at the request level (layer 7) routing traffic to targets - EC2 instances, containers and IP addresses based on the content of the request.

Ideal for advanced load balancing of HTTP and HTTPS traffic, Application Load Balancer provides advanced request routing targeted at delivery of modern application architectures, including micro-services and container-based applications.

III. PROPOSED

This paper explains how to use autoscaling to automatically adjust the number of VM instances that are hosting your application, allowing your application to adapt to varying amounts of traffic. To useautoscaling, host your application on a managed instance group. A managed instance group is a collection of instances that are all running the same application and can be managed as a single entity. When a managed instance group has autoscaling enabled, the number of VMs in the instance group automatically increase(scales out) or decrease(scales in)according to the target value that you specify for your autoscaling policy. This project includes detailed steps for launching a web application on a managed instance group, setting up autoscaling, configuring network access, and observing autoscaling by simulating load spikes and drops.

THE EXPERIMENTAL RESULTS

Fig 1: Instance Template

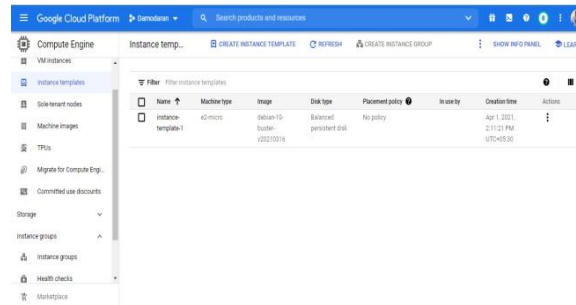


Fig 2: Instance Group

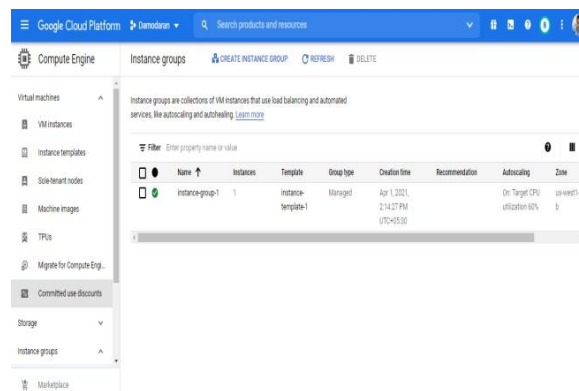
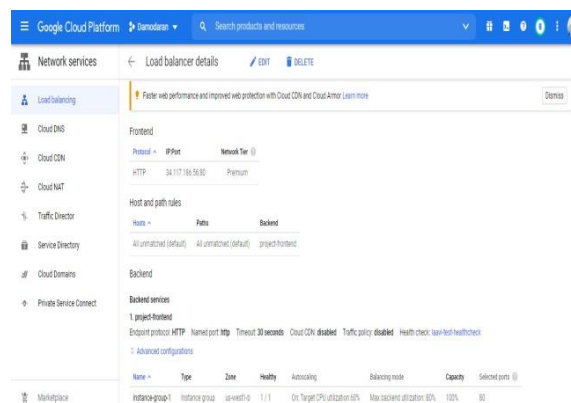


Fig 3: Load Balancing



15th & 16th September 2021

Organized by

Department of Computer Science, Parvathy's Arts and Science College, Dindigul, Tamilnadu, India

Fig 4: Healthy status

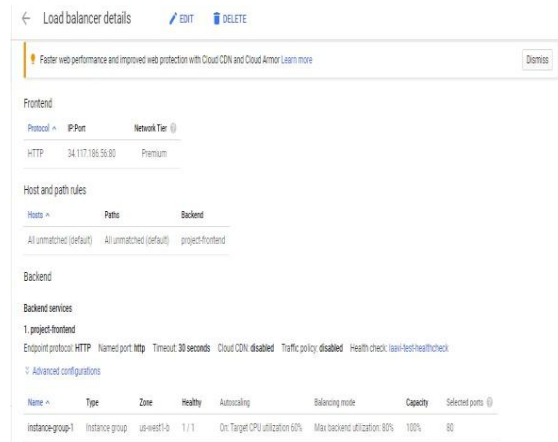


Fig 5: SSHing

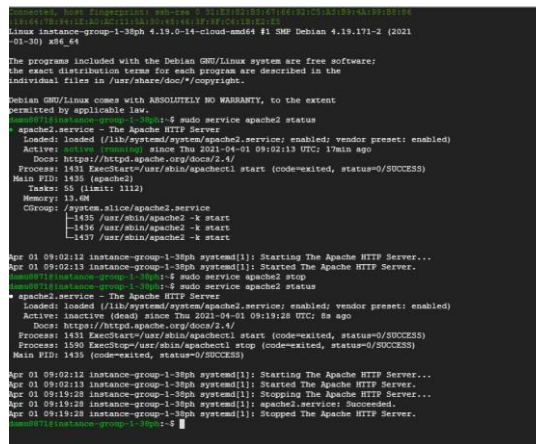
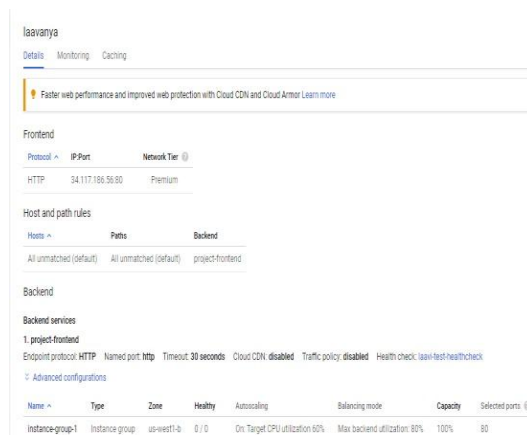


Fig 6: Unhealthy status



15th & 16th September 2021

Organized by

Department of Computer Science, Parvathy's Arts and Science College, Dindigul, Tamilnadu, India

Fig 7: Autoscale the Healthy status

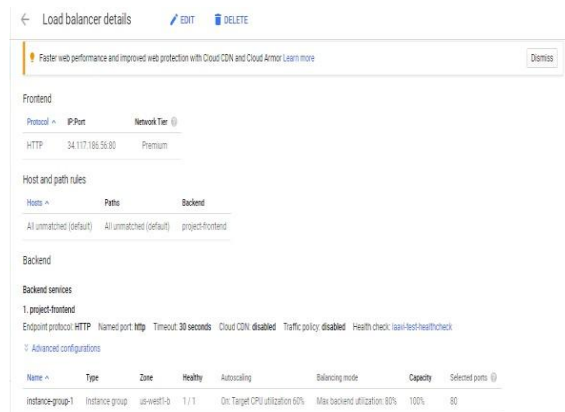


Fig 8: Cloud Shell

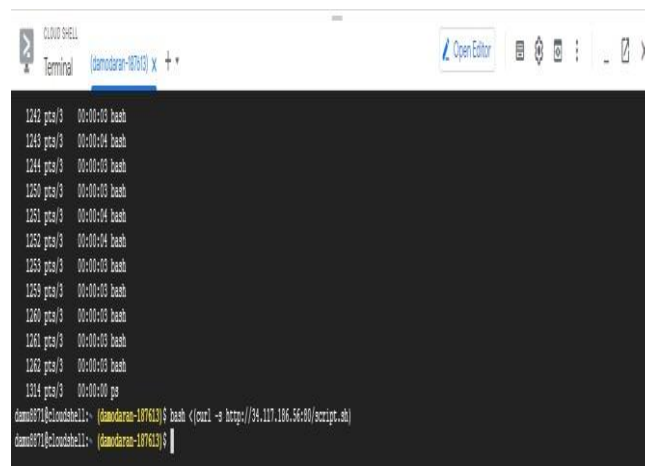


Fig 9: Autoscaling Up

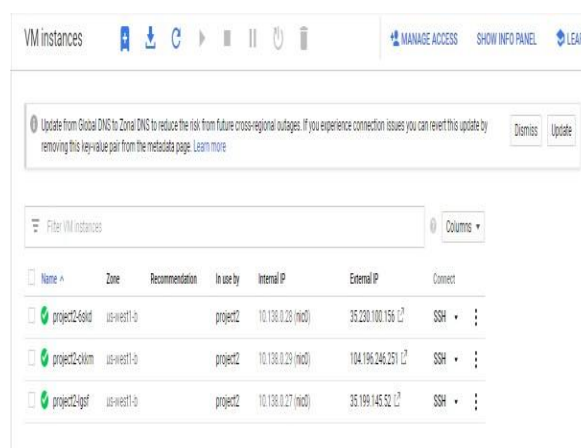
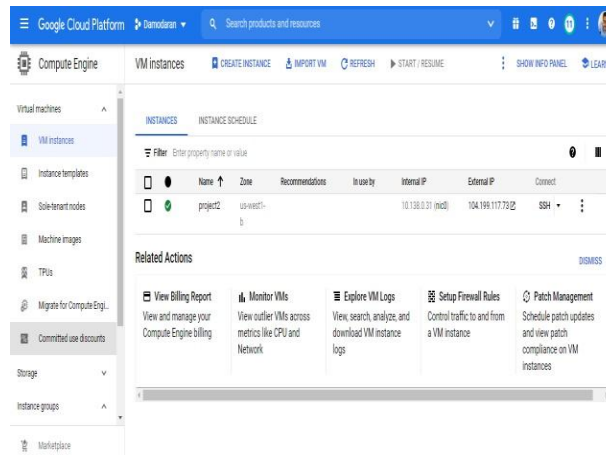


Fig 10: Autoscaling Down



IV. RESULT

The first autoscaling strategy is to simply configure the autoscaling to maintain a set number of instances indefinitely. EC2 autoscaling routinely scans instances to determine their health, if it detects a bad instance, it will end it and launch a replacement one. This gives you a predetermined number of instances, running at all time. You can always go back to manual scaling, which is a primary way of scaling resources. EC2 autoscaling can manage instance creation and termination to upkeep a stable capacity, which is a value you've specified. This allows you to maintain a maximum, minimum, or other desired capacity of your choice for your autoscaling group. Scaling events can be set to occur automatically at a certain date and time. This is especially helpful in situations where you can accurately forecast demand. What's different about this strategy is that following a schedule predicts the number of available resources at a given time in advance rather than using automation to determine appropriate mounts from moment to moment. While AWS autoscaling can perform all of the more traditional scaling methods mentioned in strategies one through three, scaling along with demand is where AWS's unique capabilities start to shine. The ability to shift seamlessly between the more traditional strategies and those discussed in number four and five is another nice feature of AWS autoscaling in and of itself. Demand-based scaling is highly responsive to fluctuating traffic and helps accommodate traffic spikes you cannot predict. That makes it a good all-around, it has a few different settings, too. For example, you can set CPU utilization to remain at 50 percent should application load shift.

V. CONCLUSION AND FUTURE ENHANCEMENT

To launch autoscaling at the most basic level, you need to determine your CPU target resource utilization level. This is the level where you want to maintain your VM and container instances. Ask yourself what percentage would you like to reach, whether its 50%, 75% or another level. Once you determine this, you then will need to create an autoscaling policy that is centered around your target utilization level. Autoscaling will then continuously monitor your MIG and collect usage information based on the policy you created. It will compare actual utilization with your target to determine which, if any, groups need to be scaled up or down to maintain the CPU utilization level as close to your specified level as possible. It does this by adding or removing VM instances from the MIG to keep your desired utilization level. You also can set up autoscaling to balance serving capacity loads that can be based on utilization or requests per seconds. You define the serving capacity of an instance through a backend service, a set of values used to connect backends and various distribution and session settings. Being able to add or remove VMs based on resource demand and traffic allows you to build a resilient, cost-effective Google Cloud infrastructure that uses just the right amount of resources at the right time for your application's workload. This intelligent, dynamic scaling tool helps you keep your actual spend compared to budget in check, reducing expenses and eliminating pay-as-you-go surprises, even when unexpected spikes occur. You can be sure you have the right number of Google Compute Engine instances available at any given time to handle an app's workload. This backend configuration to the load balancer can be enhanced to be in a user friendly way. A launch template or a programming language can be replace the manual overriding into the autoscaling group and the load balancer configuration. It can be run as an API call from any of the

machine which has the proper Identity and Access enabled. If we can concentrate on the programming side of the GCP to become a developer at the GCP Centre, we can customize our own template to run the setup quicker and easier way. The software can be published as a beta version to get the feedback from all the users across region to post their thoughts. Based on the feedback and thoughts we got from the software we can enhance our launch template or the design we are proposing to the DCP Developers team.

REFERENCES

- [1] Coutinho, E. F., de Carvalho Sousa, F. R., Rego, P. A. L., Gomes, D. G., and deSouza, J. N. (2015). Elasticity in cloud computing: a survey. *annals of telecommunications- annales des t  l  communications*, 70(7- 8):289–309.
- [2] Evangelidis, A., Parker, D., and Bahsoon, R. (2018). Performance modelling and verification of cloud-based auto-scaling policies. *Future Generation Computer Systems*, 87:629–638
- [3] Fernandez, H., Pierre, G., and Kielmann, T. (2014). Autoscaling web applications in heterogeneous cloud infrastructures. In *Cloud Engineering (IC2E), 2014 IEEE International Conference on*, pages 195–204. IEEE.
- [4] Gong, Z., Gu, X., and Wilkes, J. (2010). Press: Predictive elastic resource scaling for cloud systems. In *2010 International Conference on Network and Service Management*, pages 9–16. IEEE.
- [5] Huang, G., Wang, S., Zhang, M., Li, Y., Qian, Z., Chen, Y., and Zhang, S. (2016). Auto scaling virtual machines for web applications with queueing theory. In *2016 3rd International Conference on Systems and Informatics (ICSAI)*, pages 433–438. IEEE
- [6] Liu, X., Yuan, S.-M., Luo, G.-H., Huang, H.-Y., and Bellavista, P. (2017). Cloud resource management with turnaround time driven auto-scaling. *IEEE Access*, 5:9831– 9841.
- [7] Lorigo-Botran, T., Miguel-Alonso, J., and Lozano, J. A. (2014). A review of auto- scaling techniques for elastic applications in cloud environments. *Journal of GridComputing*, 12(4):559–592.
- [8] MLorigo-Botren, T., Miguel-Alonso, J., and Lozano, J. A. (2013). Comparison of auto- scaling techniques for cloud environments.
- [9] Mahallat, I. (2015). Astaw: Auto-scaling threshold-based approach for web application in cloud computing environment.
- [10] Manvi, S. S. and Shyam, G. K. (2014). Resource management for infrastructure as a service (iaas) in cloud computing: A survey. *Journal of Network and Computer Applications*, 41:424–440.
- [11] Mao, M. and Humphrey, M. (2011). Auto- scaling to minimize cost and meet application deadlines in cloud workflows. In *High Performance Computing, Networking, Storage and Analysis (SC), 2011 International Conference for*, pages 1–12. IEEE
- [12] Mao, M. and Humphrey, M. (2012). A performance study on the vm startup time in the cloud. In *Cloud Computing (CLOUD), 2012 IEEE 5th International Conference on*, pages 423–430. IEEE.
- [13] Wei, Y., Kudenko, D., Liu, S., Pan, L., Wu, L., and Meng, X. (2019). A reinforcement learning based auto-scaling approach for saas providers in dynamic cloud environment. *Mathematical Problems in Engineering*, 2019.
- [14] Xing, J., Zhou, H., Shen, J., Zhu, K., Wangt, Y., Wu, C., and Ruan, W. (2018). Asidps: Auto-scaling intrusion detection and prevention system for cloud. In *2018 25th International Conference on Telecommunications (ICT)*, pages 207–212. IEEE
- [15] , A. N., Son, J., Chi, Q., and Buyya, R. (2019). Elasticssc: Auto-scaling techniques for elastic service function chaining in network functions virtualization-based clouds. *Journal of Systems and Software*.
- [16] Eddy Caron: Auto-Scaling, Load Balancing and Monitoring in Commercial and Open- Source Clouds
- [17] Miss.RudraKoteswaramma : Client-Side Load Balancing and Resource Monitoring in Cloud , ISSN:2248-9622
- [18] Amazon Elastic Compute Cloud <http://aws.amazon.com/ec2/>.
- [19] Amazon web services cloud watch Web Site, November 2013.
- [20] Aws elastic load balancing Web Site, November 2013



**Impact Factor:
7.569**



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 **9940 572 462**  **6381 907 438**  **ijirset@gmail.com**



www.ijirset.com

Scan to save the contact details