



MECHANISM OF ELASTIC LOAD BALANCING

¹Mr.Dharmendra Kumar Roy,²Arigela Sathwika,³Chinta SaiDeepthi,⁴Dubbaka Poojitha,⁵Kadari Pravalika

¹Assistant Professor,²Student,³Student,⁴Student,⁵Student

¹Computer Science & Engineering,

¹Hyderabad Institute of Technology & Management, Hyderabad, India

Abstract: What if one day u wake up and your application starts getting a lot of unexpected traffic? This is amazing, isn't it? But is your application ready for it? Can it handle so much traffic? In this project we propose a Elastic load balancer. ELB (Elastic Load Balancer) is a service provided by AWS (Amazon Web Services) which is used to balance the incoming load (traffic) on multiple EC2 instances which can be in multiple availability zones. Load balancing distributes the network traffic in such a manner that increases speed and the performance of web applications. Modern ELB uses "Round robin algorithm". The real time applications of ELB are high traffic websites, E- commerce applications, Media streaming, Gaming applications, Mobile applications. ELB improves scalability, availability and fault tolerance of web applications and services.

Keywords–Elastic Load Balancer, Amazon Web Services, Elastic Compute Cloud, Round Robin, Scalability.

I. INTRODUCTION

Amazon Web Services (AWS) is a subsidiary of Amazon that provides on- demand cloud computing platforms and APIs to individuals, companies, and governments, on a metered, pay-as-you-go basis. Oftentimes, clients will use this in combination with autoscaling (a process that allows a client to use more computing in times of high application usage, and then scale down to reduce costs when there is less traffic). These cloud computing web services provide various services related to networking, compute, storage, middleware and other processing capacity, as well as software tools via AWS server farms. This frees clients from managing, scaling, and patching hardware and operating systems



Figure 1.1 Amazon web services

AWS Services

Reliability and Availability: AWS offers high availability and redundancy through its data centers, ensuring minimal downtime and improved reliability for your applications.

Security: AWS provides robust security features and tools to help you protect your data and applications. It includes encryption, identity and access management, and compliance. AWS (Amazon Web Services) offers a wide range of cloud computing services that provide numerous benefits to individuals, businesses, and organizations. Here are some reasons why people use AWS services:

Scalability: AWS allows you to scale your computing resources up or down based on demand. This flexibility is particularly useful for businesses with varying workloads.

Cost Savings: AWS operates on a pay-as-you-go model, where you only pay for the resources, you use. This eliminates the need for large upfront investments in hardware and infrastructure.

Global Reach: AWS has data centers in multiple geographic regions, enabling you to serve customers and users from around the world with low latency.

Wide Variety of Services: AWS offers a vast array of services, including computing power, storage, databases, machine learning, analytics, content delivery, and more. This allows you to choose the services that best suit your specific needs.

Innovation and Speed: AWS provides tools and services that enable rapid development and deployment of applications, allowing you to innovate and bring products to market faster.

Managed Services: AWS manages many aspects of the underlying infrastructure, such as server maintenance, security patching, and updates.

2. LITERATURE SURVTVEY

A comparative study on distributed load balancing algorithm for cloud computing.

Lua et al. [1] proposed a load balancing scheme called as Join-Idle-Queue. This scheme is based on distributed load balancing which is done by distributed dispatcher. At the 1st step all distributed dispatchers make idle processors queue. Then join these idle queues to assign the jobs coming to the servers to reduce the load of other overloaded nodes. It also claims to reduce the response time.

Chen et al. [2] proposed (1995) a load balancing scheme called as User-Priority guided min-min scheduling algorithm for load balancing in cloud computing. First, they found the minimum execution time of all tasks then the maximum value was selected. This algorithm was proposed for static environment and simulated on cloudsim.

Liu et al [3] proposed (2006) a lock-free multiprocessing load balancing scheme. This scheme reduces the use of shared memory and improves the overall performance in multi-core environment.

Nitish C et al. [4] proposed load balancing scheme HEFT based workflow scheduling for cost Optimization with in Deadline in Hybrid clouds. They have simulated their work on workflows. They have taken deadline completion as main issue in load balancing for independent tasks and tried to reduce the overall cost.

Randles et al. [5] done comparative study on distributed load balancing algorithm for cloud computing. They said that there are three methods for large scale load balancing in cloud systems- 1. Nature inspired 2. Random sampling of system domain, 3. Restructured system to optimize job assignments.

Chunling C. et al [6] proposed energy saving task scheduling strategy based on vacation queuing theory in cloud computing for dynamic environment. They have used vacation queuing model with exhaustive service to schedule the tasks. On the basis on busy period and busy time they have analyze the energy consumption of nodes. They have simulated algorithm using Matlab tool.

Brian Chiang [7] (2023) claimed Effectively balancing traffic in datacenter networks is a crucial operational goal. Most existing load balancing approaches are handcrafted to the structure of the network and/or network workloads. Thus, new load balancing strategies are required if the underlying network conditions change.

Steven Hofmeyr [8] (2010), asserted Speed balancing is also able to provide comparable or better performance than DWRR, a fair multi-processor scheduling implementation inside the Linux kernel. Furthermore, parallel application performance is often determined by the implementation of synchronization operations and speed balancing alleviates the need for tuning the implementations of such primitives.

Onur Destanoglu [9] (2008), Load balancing is performed to achieve the optimal use of the existing computational resources as much as possible whereby none of the resources remains idle while some other resources are being utilized. Balanced load distribution can be achieved by the immigration of the load from the source nodes which have surplus workload to the comparatively lightly loaded destination nodes.

Ayushi Verma [10] (2022), Cloud computing enables users to run applications and web services without having to worry about servers. However, certain occurrences, such as a power outage or a connection issue,

might cause such servers to stop functioning. Hence, it is vital to address such issues to determine the root cause of server failure and how such issues may be rectified. As a result, we employed AWS services.

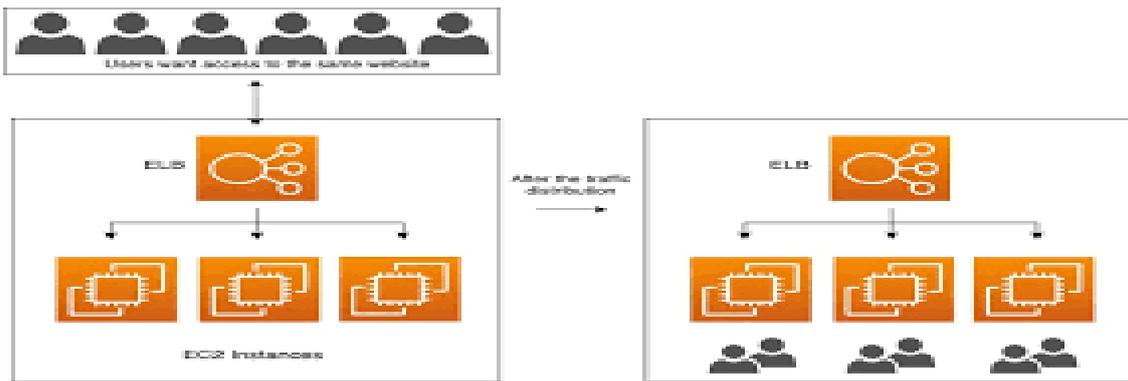


Fig 2.1 Working of ELB

In AWS, the Elastic Load Balancer (ELB) service, now known as the Application Load Balancer (ALB) and Network Load Balancer (NLB), we can use a flow-based, rather than a traditional round-robin algorithm. This flow-based approach ensures that incoming requests are directed to the appropriate target based on factors like the request's source IP, the destination IP, and the TCP/UDP port.

However, you can achieve a round-robin-like behavior in AWS ELB/ALB by configuring the load balancer's target group to use a "Least Outstanding Requests" or "Least Connections" routing algorithm. While these options don't provide a strict round-robin approach, they distribute traffic fairly evenly among healthy targets by selecting the target with the fewest outstanding requests or connections.

With the configuration of load balancer in Amazon Web Services Management Console, the ALB will distribute incoming requests to targets based on the algorithm you selected that is Round Robin which will result in a more evenly distributed load.

3. PROPOSED SYSTEM

In AWS, the Elastic Load Balancer (ELB) service, now known as the Application Load Balancer (ALB) and Network Load Balancer (NLB), we can use a flow-based, rather than a traditional round-robin algorithm. This flow-based request approach ensures that incoming requests are directed to the appropriate target based on factors like the request's source IP, the destination IP, and the TCP/UDP port.

However, you can achieve a round-robin-like behavior in AWS ELB/ALB by configuring the load balancer's target group to use a "Least Outstanding Requests" or "Least Connections" routing algorithm. While these options don't provide a strict round-robin approach, they distribute traffic fairly evenly among healthy targets by selecting the target with the fewest outstanding requests or connections.

With the configuration of load balancer in Amazon Web Services Management Console, the ALB will distribute incoming requests to targets based on the algorithm you selected that is Round Robin which will result in a more evenly distributed load.

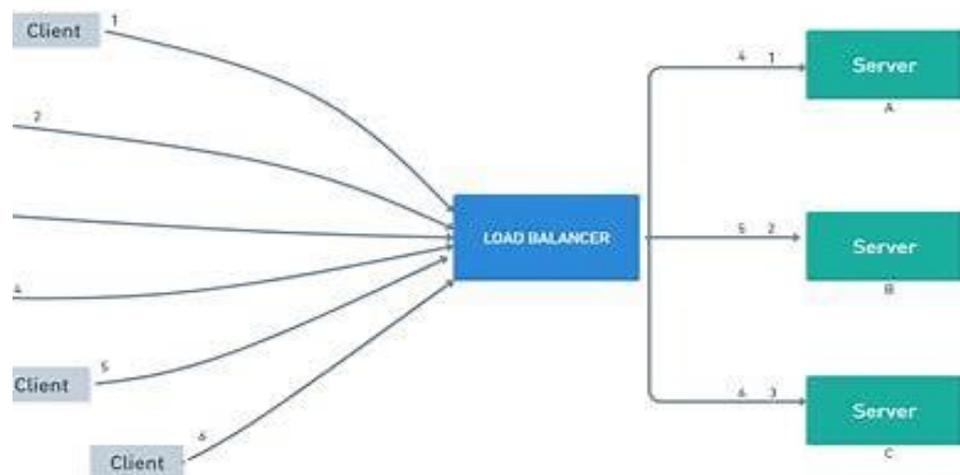


Figure 3.1 Nginx Load Balancer

Nginx:

Nginx (pronounced "engine-x") is a popular open-source web server, reverse proxy server, and load balancer. It's designed to efficiently handle HTTP, HTTPS, and other web protocols. Nginx is known for its high performance, scalability, and versatility, making it a widely used software in web hosting and server management. Nginx can act as a reverse proxy server, sitting in front of application servers (e.g., web applications, API servers) and forwarding client requests to the appropriate backend servers. This helps distribute traffic, improve security, and provide features like load balancing. It can distribute incoming traffic across multiple backend servers, ensuring that each server shares the load evenly. This is crucial for high availability and scalability in web applications.

It offers security features like access control, rate limiting, and the ability to filter and block malicious traffic. It can be used as part of a defense strategy against DDoS attacks and other security threats. Its lightweight and efficient nature make it a popular choice for serving web content and managing network traffic effectively.

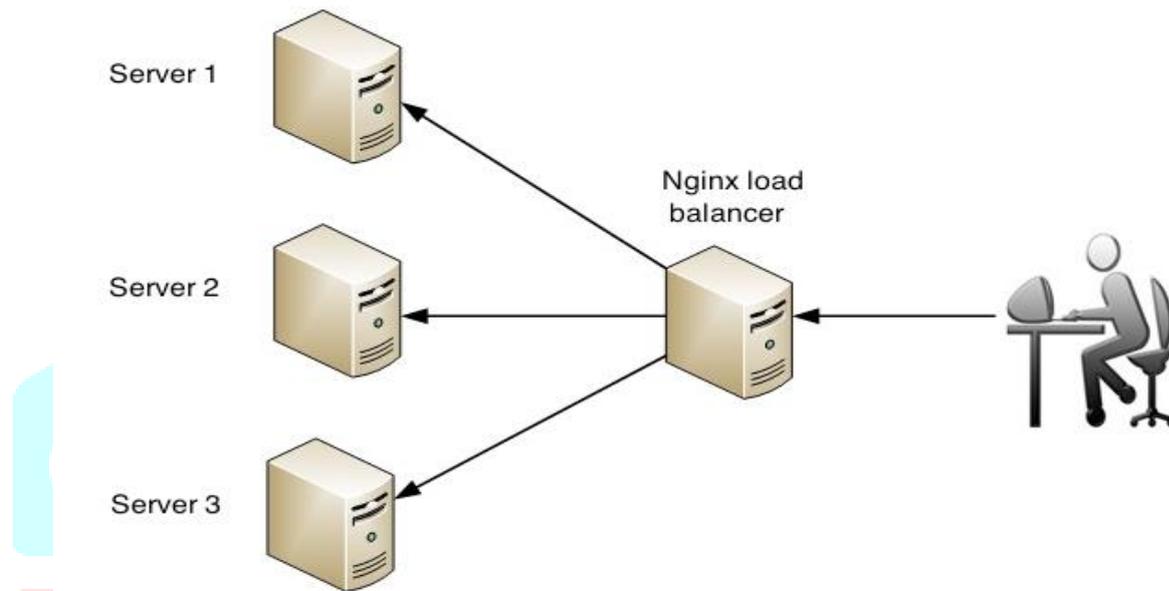


Figure 3.2 Round Robin

In documentation, a proposed approach typically refers to a recommended plan or strategy for achieving a certain goal or solving a problem. It outlines the steps, methods, and considerations that should be taken into account when implementing a particular solution. In the context of AWS documentation, a proposed approach could involve using specific AWS services, best practices, architectural considerations, and implementation guidelines to achieve a desired outcome or address a specific use case. It includes creating AWS designs and defining migration strategies and patterns that support the move to the cloud and further optimization. This type of assessment is typically time consuming, and it's approached in a sequence of smaller chunks, aligned to migrated waves, throughout the project lifecycle.

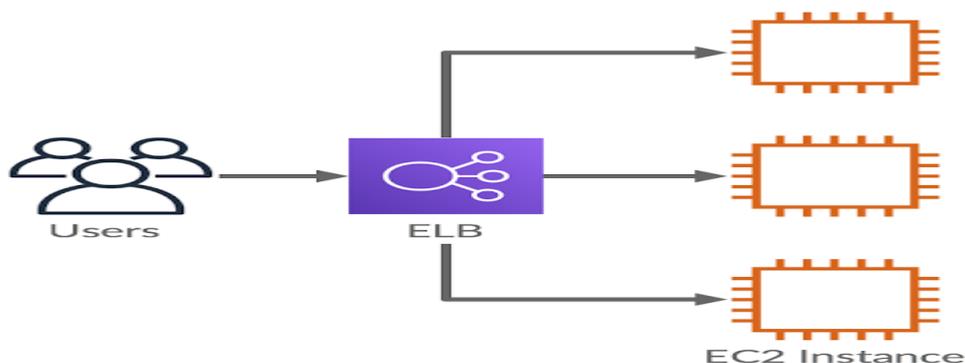


Figure 3.3 Traffic Distribution To Multiple Instance

The proposed system provides a user-friendly system through which farmers can generate detailed information. These profiles contain important information about their farm, their development and the products they offer. Farmers can list their crops, animals and other products, and provide descriptions and photos.

Existing Problem

Elastic Load Balancers (ELBs) are an integral part of distributing incoming network traffic across multiple targets, such as Amazon EC2 instances, containers, and IP addresses. One exciting problem in this context could involve optimizing the ELB configuration to achieve the best performance, high availability, and cost-effectiveness for your application. This might include:

Auto Scaling: Designing the ELB setup to dynamically adjust its capacity based on traffic spikes or drops, ensuring your application can handle varying loads efficiently.

Health Checks: Ensuring that the ELB accurately monitors the health of its registered instances and routes traffic only to healthy instances.

Cross-Zone Load Balancing: Configuring the ELB to distribute traffic evenly across multiple Availability Zones, providing better fault tolerance and reducing latency.

SSL/TLS Termination: Implementing SSL/TLS termination at the ELB to offload the decryption process from backend instances, improving performance and security.

Path-Based Routing: Utilizing path-based routing to direct different types of requests to specific backend instances or groups, enabling more efficient resource utilization.

Session Persistence: Handling session persistence effectively, either through sticky sessions or other mechanisms, to maintain user sessions across requests.

Monitoring and Logging: Setting up comprehensive monitoring and logging for the ELB to quickly identify and troubleshoot issues, helping maintain high availability.

Cost Optimization: Fine-tuning the ELB configuration to optimize costs while maintaining performance by using appropriate instance types, scaling policies, and traffic routing rules.

Geo-based Routing: Implementing geolocation-based routing to direct traffic to the nearest data center or region, improving latency for users around the world.

Customizing Load Balancing Algorithms: Exploring different load balancing algorithms offered by the ELB and selecting the one that best suits your application's requirements.

4. RESULTS

Output on Home Screen:

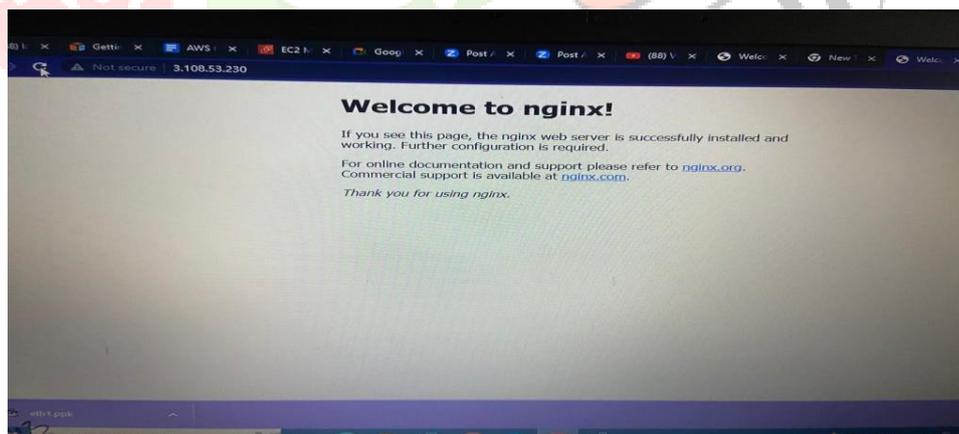


Figure 4.1: Output Without Using The URL

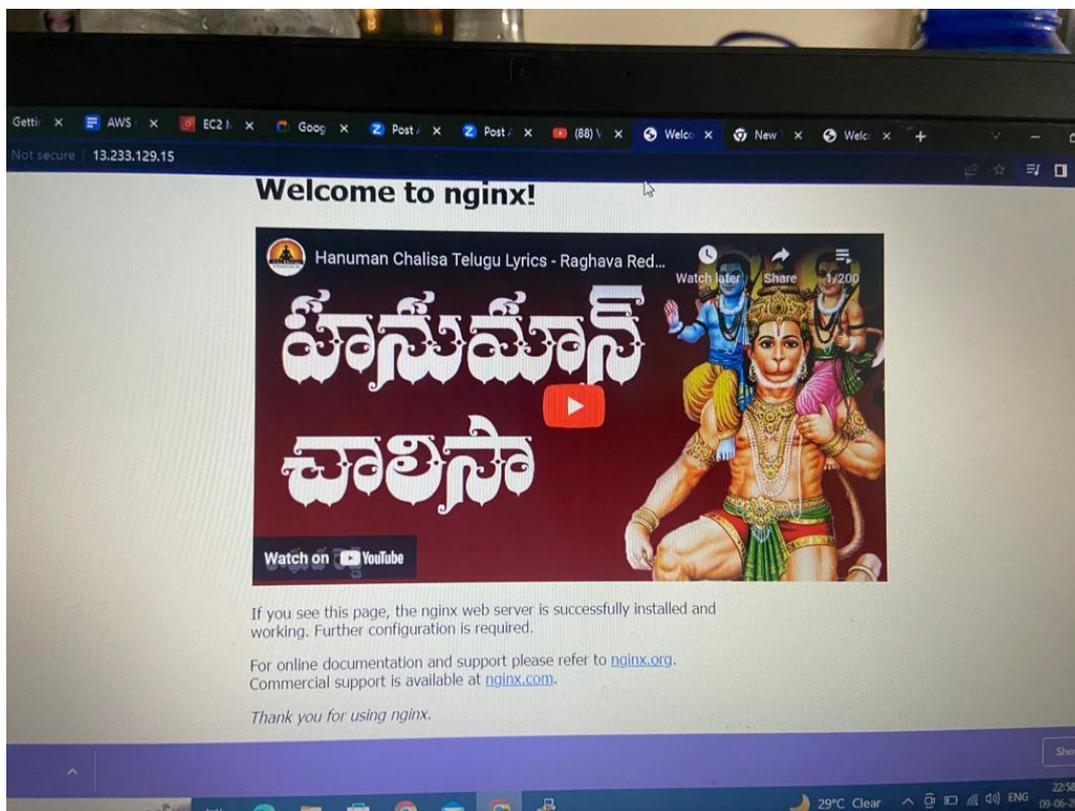


Figure 4.2 Output with the URL

CONCLUSION

In conclusion, our mini-project successfully demonstrated the implementation of an Elastic Load Balancer (ELB) in Amazon Web Services (AWS) using the round-robin algorithm, all achieved through the power of the C programming language. By leveraging AWS services and the round-robin algorithm, we've created an efficient and scalable load balancing system. Throughout this project, we've gained valuable insights into cloud computing, networking, and load balancing principles. We've learned how ELB can evenly distribute incoming requests to a group of backend servers, ensuring high availability and optimizing resource utilization. The round-robin algorithm, implemented in C, proved to be a simple yet effective means of achieving this load balancing. Implementing Elastic Load Balancing (ELB) in Amazon Web Services (AWS) with the round-robin algorithm offers a straightforward and effective approach to distributing incoming network traffic across multiple EC2 instances. This load balancing methodology ensures even traffic distribution, promoting high availability and scalability for applications. It is particularly suitable when the backend servers are homogeneous and capable of handling requests independently. By systematically routing requests to each instance in a circular fashion, round-robin load balancing optimizes resource utilization and enhances fault tolerance, contributing to a robust and efficient workload management solution within the AWS ecosystem. However, it's important to consider the specific requirements of your application and workload when choosing a load balancing algorithm to ensure it aligns with your performance and reliability objectives. While this mini-project provided a solid foundation, there is always room for further enhancements and exploration. Future iterations could involve integrating more advanced load balancing algorithms, adding monitoring and logging features, or extending the project to include other AWS services for a comprehensive cloud-based solution. In a rapidly evolving field like cloud computing, continuous learning and innovation are essential, and this mini-project serves as a stepping stone toward that ongoing journey.

REFERENCES

- [1] Lua et al (2011) "Load Balancing Algorithms in Cloud Computing Environment: A Review", International Journal on Recent and Innovation Trends in Computing and Communication, Vol 2, Aug 2011.
- [2] Chen (2014)"An Overview of Cloud Computing Adoption Challenges in the Norwegian Context." Utility and Cloud Computing (UCC), 2014.
- [3] Liu (2014) "A comparative study into distributed load balancing algorithms for cloud computing." Advanced Information Networking and Applications Workshops (WAINA), 2010 IEEE 24th International Conference on. IEEE, 2014.
- [4] Nitish (2015) "Comparative study on load balancing techniques in cloud computing." Open Journal of Mobile Computing and Cloud Computing 1.1 (2014).
- [5] M. Randle's (2010) "An in-depth analysis and study of Load balancing techniques in the cloud computing environment." Procedia Computer Science 50 (2015).
- [6] Chunling (2018)" Load balancing and resource monitoring in cloud", Proceedings of the CUBE International Information Technology Conference. ACM, 2018.
- [7] Brian Chang (2023) "Leading Load Balancing" ICDCN '23: Proceedings of the 24th International Conference on Distributed Computing and Networking January 2023.
- [8] Steven Hofmeyr (2010) "Load Balancing on Speed" PPOPP '10: Proceedings of the 15th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming January 2010.
- [9] Onur Destanoglu (2008) "Randomized Load Balancing Approach "ICPPW '08: Proceedings of the 2008 International Conference on Parallel Processing - Workshops, September 2008.
- [10] Ayushi Verma (2022) "Dynamic Target Monitoring of Load Balancers in Cloud Computing" IC3-2022: Proceedings of the 2022 Fourteenth International Conference on Contemporary Computing, August 2022.

