

E-ISSN: 2229-7677 • Website: <u>www.ijsat.org</u> • Email: editor@ijsat.org

Scaling Software Engineering in the Cloud: Leveraging AWS Auto Scaling for Dynamic Application Growth

Sai Krishna Chirumamilla

Software Development Engineer II, Dallas, Texas, USA saikrishnachirumamilla@gmail.com

Abstract:

Over the last few years, cloud computing has revolutionized how software engineering is done by providing elastic, affordable, and flexible resources. This paper focuses on AWS Auto Scaling as an enabler of dynamic growth, resource optimization, and creating fault tolerance for the software engineering teams. Our systematic literature study discusses patterns, issues, and recommendations for auto-scaling approaches. Theoretically, the paper is methodologically focused on illustrating a real-life application of AWS Auto Scaling deployment in detail. Several outcomes show that such an approach helps enhance response times, minimize downtimes, and maximize costs. These results shed light on the cloud technologies' prospect for engaging emergent concerns in the modern software development sphere.

Keywords: AWS Auto Scaling, Cloud computing, Dynamic scaling, Software engineering, Fault tolerance.

1. Introduction

1.1. Evolution of Scaling Software Engineering in the Cloud

The transition to scale software engineering in the cloud has been one of the most exciting transformations in application development, deployment, and management. [1-4] Cloud technologies as products evolution path from monolith architecture up to today's cloud-native approaches have significantly enhanced the method and manner of resources/applications scalability levels that has never been possible within traditional on-premises approaches. The following sections explore key stages of this evolution:



E-ISSN: 2229-7677 • Website: www.ijsat.org • Email: editor@ijsat.org



Figure 1: Evolution of Scaling Software Engineering in the Cloud

- The Emergence of Cloud Computing (Pre-2000s): However, with the use of the cloud computing model, scaling of the software application was very manual and needed the procurement of physical infrastructure. Applications would require organizations to acquire and maintain dedicated hardware for their applications, which was costly and time-consuming. With the advent of cloud computing, earlier versions, such as AWS and Microsoft Azure, changed the traditional computing model by presenting a flexible model for computing and storage services on the cloud. However, scaling was still primarily static at this level, and resources were mainly preconfigured or allocated in large blocks at one time rather than in response to the current workload. These first cloud solutions were revolutionary in terms of the way that IT infrastructure was consumed and delivered and provided a base upon which future application architectures could be built.
- Monolithic to Modular Architectures (Early 2000s): It was around the year 2000, and as cloud computing started to take a foothold, software development also migrated away from more centralized approaches to monolithic software development. For example, applications with tightly interlinked components such as UI, business logic, and databases in monolithic systems were difficult to scale. When the demand rose, full systems had to be copied, which was time-consuming and costly. Solutions for constructing applications on private or public clouds enabled developers to split an application into small, autonomous services or microservices that can be scaled apart. These changes continued the process of modernization, which solved some of the problems of monolithic systems, providing more specific scaling that only created copies of necessary parts of an application, avoiding overheads and increasing flexibility.



E-ISSN: 2229-7677 • Website: www.ijsat.org • Email: editor@ijsat.org

- Introduction of Virtualization and Containerization (The mid-2000s): The next explosion in cloud scaling came with the advent of virtualization technologies in the mid-2000s. It enabled multiple virtual machines operational in a solitary physical host; however, they are logically severed all through the process of virtualization, multiple. This was becoming beneficial in allowing developers and the IT departments to offer dynamic solutions for scaling the applications. This was later done away with by containerization as another method of changing this approach. Docker, launched in 2013, enables the implementation of the application and many of its dependencies in the form of lightweight containers whose homogeneous environments can be launched on different platforms. Containers are lighter than VMs because they run on a similar operating system kernel; container-based virtualization is also the most efficient, taking less time to deploy. This laid the foundation for the next step in cloud scaling: management, provisioning, and automatic scaling of applications that become containers.
- The Advent of Microservices and Distributed Systems (Late 2000s 2010s): In the late 2000s and 2010s, a new approach to scaling applications, known as microservices architecture, was introduced. Instead of scaling whole complex monolithic systems, the organizational behavior shifted to decomposing apps into numerous separate and standalone services, each of which handles one microservice. This made it easy to bring scaling for each microservice up or down depending on the traffic received and the overall demand. Thus, Kubernetes, as shown below, or Amazon ECS (Elastic Container Service), was developed to help manage these distributed systems' self-scaling containerized services. An important advantage of microservices was that it was easy to control the scaling because each service could be independently scaled, while the cloud platforms granted the application the ability to scale this single service without scaling the rest of the application.
- Cloud-Native Applications and Serverless Computing (2010s Present): Eligibility was first observed in the 2010s as a feature of software developed to be native for the cloud from the ground up, able to run efficiently and scale on the cloud. Cloud-native applications make use of microservices and containers; many cloud-native applications are constructed based on CI/CD principles. Perhaps one of the biggest revolutions in this period was the idea of serverless, which was made famous by services such as AWS Lambda (2014) that let code be executed without the need to know or manage servers. In a serverless model, paired with the idea of the application's resources being provided by the cloud provider, the applications, in turn, are self-scaling. This approach is also inherently efficient for organizations' finances because organizations are charged for specific compute time and not for unused server space. Serverless computing has done a very good job of creating an environment for applications to scale because they can scale independently without human intervention.
- The Role of Auto-Scaling in Cloud Computing (2010s Present): When cloud computing evolved over the years, auto-scaling was to be part of most of the cloud technologies. Auto-scaling entails the application of resource capacity, including virtual machines and containers at specific utilization as per the real-time metrics of computations involving the CPU, RAM, and network traffic. For instance, AWS Auto Scaling adjusts the quantity of EC2 instances deployed horizontally through addition or removal depending on the workload. These make scaling less about manually charting courses and more about an efficient and responsive process. This makes resource utilization efficient and considerate of operating expenses since Auto-scaling only provides resources when required. By integrating auto-scaling into the cloud services, organizations have been in a better position to



address multi-tenanted issues with traffic variations while at the same time increasing service performance and efficiency.

• Challenges and the Future of Cloud Scaling: Still, several areas need improvement, even under the context of massive improvements in cloud scaling. Being able to scale up has been a concern with many systems where the actions are slow and at times when the new instances or containers need to be provisioned to handle the increasing demand. These delays can lead to performance degradation, at least temporarily, or service interruption. Also, designing proper scaling policies to achieve the right scales between response time and cost may be challenging because it depends on traffic behavior, resource usage, and SLAs. In the future, scaling accuracy could benefit from the ML and predictive analytics concepts, especially through integrating both. Based on past trends data analyses by the ML models, more resources could be procured before a spike in demand was foreseen, thus helping avert timing lags, which otherwise would drag down system efficiency. Moreover, the expansion of the requirement for using multiple clouds and the rationalization for the further development of the serverless computing concept can indicate the emergence of sophisticated contextual scaling solutions for various workloads and configurations of environments.

1.2. Importance of Leveraging AWS Auto Scaling

AWS Auto Scaling is an enhancement that allows organizations to automatically add or remove instances such as EC2 instances, databases, or containers. [5,6] AWS Auto Scaling offers organizations the benefits of achieving application performance, effective utilization of resources, and affordable expenses due to the ability of automatic scaling. The importance of leveraging AWS Auto Scaling can be discussed across several key aspects:



Figure 2: Importance of Leveraging AWS Auto Scaling

- **Improved Application Performance:** AWS Auto Scaling even boosts the power of the application since it provides and tears down resources on the spot, depending on the load. During hours of high traffic, say during sales or promotion, Auto Scaling easily adds a computing resource so that applications can easily handle the traffic without slowing down. This dynamic scaling prevents infrastructure overloading, for example, slow page loads or even a site going down in front of the users. Further, Auto Scaling scales down during the less intense periods so that the application will not have to eat up the expenses it does not need continually, but it has a seamless run.
- **Cost Efficiency and Optimization:** AWS Auto Scaling helps to determine resource usage patterns and provides flexibility in cost control rather than predefining the resources needed. This means businesses can now avoid committing to reserving for peak demands since those result in procuring costly over-provisioned infrastructure that is not fully utilized. However, through AWS Auto



E-ISSN: 2229-7677 • Website: <u>www.ijsat.org</u> • Email: editor@ijsat.org

Scaling, what is allowed is the pay-as-you-go modality; that is, EC2 instances should be scaled up during periods of high traffic and scaled down during times of low traffic. This elasticity also assists organizations in reducing wastage instances and only billing the consumer based on the amount of usage, hence cutting on their expenses and increasing usage of resources tremendously.

- Increased Availability and Reliability: AWS Auto Scaling enhances applications' availability and reliability by ensuring that applications can always access adequate resources to serve traffic spikes. The second key feature of Auto Scaling is making traffic distribution of incoming requests across multiple instances protected by AWS Elastic Load Balancing (ELB), thus enhancing fault tolerance and preventing one instance from becoming a bottleneck. Auto Scaling possesses a mechanism in which a failed instance is instantly replaced with a healthy instance; thus, the downtime is as low as possible, and the service is continuously available. This makes the system robust, and application availability is guaranteed during difficult situations.
- Flexibility and Adaptability: AWS Auto Scaling is very flexible since organizations can only scale resources in areas where necessary. This could be a web application, a data processing task, or a backend service. Auto scaling can bring resources depending on different performance characteristics, including CPU usage, memory usage, or any CloudWatch custom metrics of any application. It also offers both vertical and horizontal scaling that allows adding or removing instances and ensures high scalability of the application. This means AWS Auto Scaling can be used for applications that require an incline scale of resources and applications that require large amounts of resources with fine-tuning of scale.
- Seamless Handling of Traffic Spikes: Routinely, AWS Auto Scaling is more beneficial in applications facing irregular or burst traffic rates. For instance, many e-commerce sites have noticed that traffic rates are much higher during sale seasons or holidays. Auto Scaling comes in handy as HIPAA's resource capacity scales up to address these spikes without intervention to ensure the application can handle the extra load. Amazon Web Services Auto Scaling guarantees that more webpages, solutions, and apps will be capable of handling more visitors, requests, and traffic because resources for development like EC2 instances or containers are automatically ramped up when sites and services get bottlenecks or become slow and erroneous, or unresponsive.
- Simplified Resource Management: Coordinating several scaling operations might be cumbersome and sometimes ineffective, especially if the processes are manual and across a large-scale application. AWS gives this an easy way out with the help of auto-scaling, which means that the respective resources are adjusted to certain policies. Auto Scaling is the capability to specify policies that instantly initiate scaling procedures when specific metrics, like CPU or memory usage, exceed given limits without dealing with every server separately. This automation not only saves time but also minimizes the problem of human error, thus allowing the engineering teams to do more critical and creative work while optimizing the resources being used.
- **Better User Experience:** Used to sustain the right number of resources across all possible circumstances, AWS Auto Scaling balances optimal response to user input. People use the application, web service, or backend system and anticipate high performance, whether a mobile application or a service. Auto Scaling ensures resources' scale according to access demand, preventing slow response time, timeouts, or system crashes. This helps customers Change their inputs into our applications without interruption, thus improving their satisfaction and making them



stay longer. A good, uninterrupted user experience results in higher customer trust and engagement, hence better business results.

- **Real-Time Monitoring and Metrics:** AWS Auto Scaling is rather tightly connected with AWS CloudWatch, which gives detailed real-time metrics for application and resource usage. AWS CloudWatch monitors metrics such as CPU usage, memory usage, disk read/write, and network traffic, which automatically scales up or down based on certain limits. This real-time usage also ensures that resources are always available in the current usage, while applications can always operate optimally. Another type of CloudWatch log is useful for pinpointing performance issues, studying scaling patterns, and improving scaling rules for better resource utilization to make better future choices.
- **Support for Hybrid and Multi-Cloud Architectures:** The AWS Auto Scaling is not only for scaling within AWS but for Hybrid and Multi-cloud environments where organizations can scale resources over On-premises plus AWS plus other cloud. This is particularly advantageous for organizations migrating to the cloud or using systems across several platforms. By using AWS Auto Scaling, the resources required for various environments are well-sized to achieve gain and effectiveness. This characteristic of microservices makes it easier to scale over a hybrid or multiple clouds, and its flexibility and scalability can cater to different workloads.
- Enhanced Security and Compliance: Thus, the general key benefits of AWS Auto Scaling are not only associated with scaling resources to correspond to demands for resources while consuming fewer resources than a client but also with improving security and compliance since it automatically adjusts resource levels of and, thus, minimizes the attack surface. This way, it keeps the number of targets for possible hacking or exploitation to the least possible with only the number of instances and services needed. Auto Scaling also includes security features, and it works with other AWS security features like AWS Identity and Access Management (IAM) and AWS Key Management Service (KMS). Moreover, Auto Scaling offers informative logging and the possibility of auditing, helping organizations meet regulations under GDPR, HIPAA, and SOC 2 and keep resources' utilization safe and compliant.

2. Literature Survey

2.1. Evolution of Cloud-Based Software Engineering

The availability of cloud computing has brought about a shift in most software engineering practices. Software development used to be in monolithic models in which all parts are interconnected and exist within a single server. However, with the advent of cloud services, the development has moved towards microservices and serverless architectures. [7-11] Cloud computing brings incredible flexibility and elasticity because firms can quickly develop and expand their applications based on customer needs. Implementing software applications as distributed systems on the cloud allows software engineers to develop applications and let the components work in parallel. This makes them modular, independently deployable, and scalable – thus more flexible and fault-tolerant. This shift has been instrumental in allowing organizations to adopt DevOps practices and embrace practices that support continuous delivery of an improved pace and quality of software deployment. With new changes in cloud services, the flexibility factor is further propelling the usage of containerization, orchestration, platforms such as Kubernetes, and serverless functions for software development and operations.



E-ISSN: 2229-7677 • Website: <u>www.ijsat.org</u> • Email: editor@ijsat.org

2.2. Dynamic Resource Scaling

Resource dynamism has emerged as one of the fundamental models of cloud systems. One of the most compelling advantages of the cloud solution is its ability to manage resources according to real users' needs at any given time without requiring additional hardware investments. Systems with auto-scaling capabilities are far more highly effective than those without regarding their availability and efficiency of their usage. Their study shows that service-based applications with dynamic scaling attributes work up to 40 percent more reliably than statically scaled applications incapable of accommodating traffic fluctuations. Auto-scaling defines the best way to provide resources, avoiding scenarios where more resources are provided than necessary and avoiding very high costs. The research also expresses how the scale-up and scale-down features in the cloud reduce resource utilization costs during low demand. In many applications, dynamic scaling has emerged as an important characteristic of modern applications as more companies embrace cloud infrastructure for running their digital operations.

2.3. AWS Auto Scaling in Practice

AWS Auto Scaling is one of the services launched in 2018 and allows users to scale EC2 instances, Amazon RDS databases, and containers in response to demand. Stated that several research studies have claimed that AWS Auto Scaling helps a great deal in meeting Service Level Agreements (SLAs), especially during traffic and usage bursts. The empirical studies presented in their research show that AWS Auto Scaling is useful for making applications always available, timely responding to incoming traffic and does not affect response time when traffic increases. The AWS Auto Scaling enables scalability through horizontal scalability, in which the system changes the instances of resources being used. It is more useful in applications that may have highly fluctuating traffic since it ensures that the resource capacity is always set to handle the load so that there will be no more than one bottleneck. AWS Auto Scaling is critical in applications that cannot afford to go down because the cost is very high, while onfocus uptime performance is critical for user satisfaction.

2.4. Challenges in Auto Scaling Implementation

However, like any other solution, the application of auto-scaling, including the AWS Auto Scaling solution, has its drawbacks. The first one is the concern with the slowness of scale, including a susceptibility to a slow response to increase or decrease the number of resources devoted to a particular project due to demand shifts. This can lead to a temporarily reduced throughput that the system experiences when the necessary resources for the new level of load have not been provided. An additional limitation is that it is not easy to set the scaling policies at work as they can be intricate in many cases. Auto-scaling is sensitive to parameters, which include CPU capacity, memory usage, and scaling points, among other parameters. The inadmissible configuration has adverse effects that often result in overscaling, which means using resources that one does not need in the first place, and underscaling, which means using where inadequate resources lead to performance bottlenecks and system slowdowns. Through researching the current trends in the industries, most organizations struggle to find a balance when it comes to implementing policies in these systems to meet the high responsiveness while at the same time being cost-effective. Preliminary scaling policies, including scaling actions prior to being triggered by traffic volume and statistical analysis of previous traffic patterns, have been proposed as some remedies. These solutions enable the 'right-sizing' of actions with a better probability of the system being primed for demand before a peak response is experienced, thus mitigating latency and optimizing resources.



- 3. Methodology
- 3.1. System Architecture



Figure 3: System Architecture

- Frontend: Hosted on EC2 Instances: The frontend layer of the application is managed for communication with the end-users and for display of the UI. It is run on Amazon EC2 instances to take advantage of the flexibility, horizontal scalability, and level of control provided in AWS. [12-16] These resources, EC2 instances, are programmed optimally to handle HTTP requests and ensure that user experience is not compromised. This way, the front end can always be running as part of an Auto Scaling group so that load balancing across multiple instances would provide low latency and high availability. Also, Amazon Route 53 is used as domain management, which helps users route to the nearest host EC2 instance.
- **Backend: Managed via an Elastic Load Balancer (ELB):** This backend contains all the business logic and processes the user's request sent from the front end of the application. A load balancer of the elastic variety is used to spread the incoming traffic across the servers on the back end. This helps prevent loading path criteria that could be detrimental to server loads and overall reliability. Second, ELB also helps in fault tolerance by providing an instance for redirecting traffic in unhealthy instances. When used with the Auto Scaling group, the ELB strives to retain performance during periods of high volume by guaranteeing new instances are registered with the load balancer.
- Database: Amazon RDS for Scalable Data Storage: At this layer, the storage and retrieval of data for the application are done and can be done in any database management system. This is made possible by employing Amazon Relational Database Service (RDS), which offers such benefits as scalability and reliability, apart from easy management. RDS supports several database engines, and MySQL is used for structured data processing in this architecture. To achieve high availability, the RDS instance has been set up as Multi-AZ, which can failover to another readily available instance. Backup and snapshots also contribute toward improving data durability and decreasing the management overhead. The authors also assert that to scale the database, they use read replicas to handle more read operations during periods of congestion.

3.2. AWS Auto Scaling Configuration

• Auto Scaling Groups: Configured for EC2 Instances: AWS Auto Scaling has a flexible workhorse called Auto Scaling Group (ASG), which lets applications keep performing optimally by automatically adjusting the number of EC2 instances in the ASG in response to the current traffic. For this case, we formed an ASG in order to have a multitude of EC2 instances for running the application frontend and backend. To maintain optimal proportional resource utilization across



different workload intensities, the optimal number of predefined ASG instances is set as minimum, maximum, and desired capacity. The ASG performs health checks on instances on a continuous basis and adjusts for groups whose instances are unhealthy by terminating and replacing them.



Figure 4: AWS Auto Scaling Configuration

- Scaling Policies: Utilized Target Tracking and Scheduled Scaling Policies: To ensure seamless resource management, the ASG is equipped with two types of scaling policies: target tracking policies and, on the other, the scheduled scaling policies. Target tracking policies automatically get the number of instances in a particular performance target; for example, average CPU utilization below 70% is achieved. These fluctuations allow the application to effectively manage different types of workloads. Scheduled scaling policies, on the other hand, are defined ready beforehand based on certain traffic patterns, like bringing up capacity for business while it is less busy, like during the night. These policies assist in avoiding compromising on either the responsiveness of the services or the cost of providing those services.
- CloudWatch Alarms: Monitored Metrics Such as CPU Utilization and Memory Usage: AWS CloudWatch also has a critical use in Auto Scaling in that it offers real-time performance data and options for the system. CloudWatch Alarms are set to monitor such important factors as the CPU percentage, the space used by memory, and the request latency. For example, if the CPU utilization is above 80%, there is an alarm that signals the launch of new instances in the ASG. On the other hand, when the utilization goes below this predetermined value, the alarm instructs the ASG to stop unnecessary instances. Such alarms allow for better timing of the scaling actions as well as for some specific requirements of application costs of latency to be effectively kept to minimal levels.

3.3. Implementation Steps



Figure 5: Implementation Steps

• Create an EC2 Instance Template with Predefined Configurations: The first requirement for using AWS Auto Scaling is to develop a reusable EC2 instance template that defines the size of the instance, OS, storage choices, and network options. These launch templates guarantee some form of standardization throughout the occurrences in the Auto Scaling group. For this setup, an AMI that has been set up with application software and necessary dependencies is employed. That is why the template also contains security group settings for the management of inbound and outbound traffic in



E-ISSN: 2229-7677 • Website: www.ijsat.org • Email: editor@ijsat.org

accordance with organizational security policies. In simplifying configurations, the template decouples the deployment process and minimizes any probability of wrong configurations.

- Set up an Elastic Load Balancer (ELB) to Distribute Traffic: After that, an Elastic Load Balancer is implemented in order to ensure that the incoming traffic is balanced between the EC2 instances forming the ASG. ELB helps make the application a single entry point, and it always routes user requests to healthy instances. HTTP, HTTPS, and TCP are among the supported protocols that explain why they are flexible enough to meet different applications' needs. In this setup, you have an Application Load Balancer (ALB), as it provides more flexibility in terms of routing requests with URLs. The ELB is also equipped with health checks to test instance health amid constant monitoring and to vote out non-responder instances from traffic.
- **Define an Auto Scaling Group with Scaling Policies:** Next, the load balancer is added, and then an ASG is configured to control the scaling activity. Finally, the ASG connects with the Ec2 instance template and the ELB for easy scaling. The minimum, maximum, and desired capacity parameters are set up for the actual management of resources. Scaling policies are then implemented: target tracking scaling policies change the number of target groups to meet real-time needs, such as CPU usage, while the scheduled scaling policies increase or decrease capacity at a given time. These policies mean that the application can deal with varying workloads as well as costs at the same time.
- Monitor Performance Using AWS CloudWatch: Last, the AWS CloudWatch is then set up to watch over the performance of the application and come up with some recommendations. The metrics, including CPU, memory, latency, and health of instances, are reflected on customized tiles in the dashboard. CloudWatch Alarms are set to start scaling actions for the pre-specified threshold or limit that is overpowered by traffic and launching Termination in the case of low traffic. CloudWatch uses logs and reports to provide system information needed to address performance issues in terms of bottleneck and scaling strategies.

3.4. Metrics for Evaluation



Figure 6: Metrics for Evaluation

• **Response Time: Measured Using Apache JMeter:** Response time is a good measure to use as the amount of work that needs to be done plays out in differing capacities. In this study, the Apache JMeter tool is used to emulate user traffic on the application and determine the time it takes for the application to process requests. JMeter creates load tests, which copy the usual genuine situations and common user traffic flows, such as many users simultaneously using the application during business hours. These tests assist in determining the behavior of AWS Auto Scaling when it comes to changing resources to respond to high traffic with low response time. Small values of response



E-ISSN: 2229-7677 • Website: www.ijsat.org • Email: editor@ijsat.org

time equal increased scalability and improved user experience, while high values may denote possible bottlenecks or the need for policy change.

- Uptime: Calculated from CloudWatch Logs: Uptime refers to how available the application is over a given period and is used to determine the calendar or company time that the application is functioning and can be used to build user confidence as well as reviewing the companies' Service Level Agreements (SLAs). This metric is pulled from the AWS CloudWatch logs, where all the health statuses of each of the EC2 instances within the Auto Scaling group are logged. Considering these logs, the uptime which the application spent in the online mode is found out. A higher percentage of uptimes indicates that Auto Scaling works efficiently to keep systems and instances failure-tolerant and frequently replaces unhealthy instances. This metric is most important for applications with high availability requirements.
- **Cost Savings: Derived from AWS Billing Reports:** Reduced costs are a quantitative measurement of the economic rationality of the application of AWS Auto Scaling. This metric is calculated by comparing the consolidation of AWS billing reports monthly pre- and post-Auto Scaling configuration. The savings are in a fraction range from the de-provisioning of excess resources at periods of low utilization. For instance, while operating during low traffic volume, the ability to set an instance down from two hundred and fifty to fifty decreases the expenses for operation without affecting the performance of an application. It has been observed that cost considerations are inherent in avoiding over-committing resources, making auto-scaling financially viable for variable workloads.

4. Results and Discussion

4.1. Performance Analysis

This means that speculative, the use of AWS Auto Scaling turned out to be significant as it improved multiple core indicators that defined the capacity of the system to work with dynamic loads. Below is an in-depth analysis of each metric:

Metric	Before Scaling	After Scaling	Improvement (%)		
Response Time (ms)	450	120	73%		
Uptime (%)	92	99.8	7.8%		
Cost (\$/month)	1500	1100	26.6%		

I ADIC I. I CHUI MANUC ANALYSIS	Table	1:	Performance	Analysis
---------------------------------	-------	----	-------------	----------



Figure 7: Graph representing Performance Analysis



E-ISSN: 2229-7677 • Website: www.ijsat.org • Email: editor@ijsat.org

- **Response Time:** In particular, with no use of the AWS Auto Scaling feature, the average response time has been determined at 450 ms, which means that the application might take time to process users' requests during peak usage. Such latency was mainly due to inadequate resources because the static environment failed to respond to changes in demand. The post-implementation response time was approximately 120 ms, which was a 73% improvement compared to the baseline. This reduction in latency was possible because the Auto Scaling group ensures that more EC2 instances are available during traffic to ensure an increase in computing power to accomplish user requests.
- Uptime: Another key area of improvement was the system uptime of the application, which records the availability of the application in question, which went up from 92% to 99.8% an improvement of 7.8%. Before using Auto Scaling, the application shut down occasionally during high load and, of course, stopped working when some instances were not replaced immediately. In Auto Scaling, the system was nearly always available because it routinely created new instances in place of unhealthy ones. This improvement was especially seen under conditions of increased traffic, where smooth vertical scaling averted server saturation and frequent downtimes.
- **Cost Efficiency:** Savings on cost remained as one of the gains from using AWS Auto Scaling. Monthly operating costs that earlier stood at \$1500 fell to \$1100, a decrease of 26.6%. This reduction was made possible by the flexibility offered by the Auto Scaling group, which was able to adjust the resources down during a less busy business period and up during a busier business period. As opposed to the pre-scaling environment where resources needed to be held in a constant number, regardless of the load, Auto Scaling guaranteed resource utilization efficiency, thus minimizing wastage.

Tuble 2. Cost Dicardown Comparison					
Metric	Before Scaling	After Scaling	Savings (%)		
Peak Hour Cost (\$/day)	50	35	30%		
Off-Peak Hour Cost (\$/day)	20	10	50%		

Table 2. Cost Breakdown Comparison



Figure 8: Graph representing Cost Efficiency

4.2. Cost Optimization

In particular, the use of AWS Auto Scaling was concluded to be highly successful due to its ability to match resource consumption to needs that exist in the actual operational environment in order to manage necessary costs. Before Auto Scaling, for static infrastructure, a certain number of EC2 instances had to be maintained to manage the traffic during their high traffic, and in the meantime, many other resources had to be wasted during low usage. Auto Scaling made the system able to downscale during low traffic



periods such as 11 PM when traffic density would reduce by a factor of twenty. During these periods, instance capacity was set lower by about 60 percent, thereby greatly improving cost efficiency without significant impacts on application throughput.

Table 3: Cost Optimization

•		
Time Period	Instances Active	Cost (\$)
Peak (8 AM - 6 PM)	5	750
Off-Peak (6 PM - 8 AM)	2	350
Total (Month)	-	1100





For instance, while handling high traffic as might be experienced during working days from 8 am to 6 pm, the number was set at five EC2 instances; on the other hand, low traffic that is witnessed at 6 pm to 8 am was comfortably handled by only two instances. This reduction simply ensured the optimization of the use of all resources in terms of compute hours, storage, and network expenses, to name but a few. Making dynamic changes in operation costs for a month reduced operational expenses from \$1500 to \$1100, meaning a 26.6% cut. This not only helps to control the amount of resources wastage but also allows the organization's control over budget savings to be extended to essential aspects like performances and disaster recoveries. It also shows the financial factors revolving around the elasticity offered by AWS Auto Scaling as well as the overall efficiency.

4.3. Challenges Encountered

• Scaling Delays: During the implementation of AWS Auto Scaling, one of the issues realized was that the scaling actions took between 2 and 5 minutes to happen. This delay primarily stemmed from two factors: boot-up time and the integral connection with the elastic load balancer (ELB). At one point, when the system had recognized a rise in traffic and more instances were required for distribution, starting new EC2 instances in Amazon Web Services required a few minutes to come



E-ISSN: 2229-7677 • Website: <u>www.ijsat.org</u> • Email: editor@ijsat.org

fully online. Also, the load balancer took some time to detect the new instances and direct traffic to them, as well. While using the application during this window, the performances of the application were slowed most especially during onslaughts of traffic when the scaling action had not come into force. While the system proved capable of accommodating the expanded usage, this lag could result in displayed user lag or failed requests during periods of high activity.

• **Complex Configuration:** Fine-tuning the scaling policies for AWS Auto Scaling, which lies within the effective configuration of the service, was also an interesting problem for optimization because of the ever-present best-of-everything dilemma of being highly responsive to application traffic changes while at the same time maintaining reasonable costs. Firstly, the appropriate value limit for the like number of CPU utilization and memory was found after testing with a specific sequence of dynamics of traffic load since their patterns were high-variant. For example, it is possible to set the scaling trigger to low, which in turn would result in unnecessary scaling up or scaling down, which is costly or, conversely, setting it high means that it may not be able to supply the necessary resources during peak usage leading to poor performance. This meant periodically throughout the day referring to traffic history logs and then reapplying the scaling policies from there in a more scaled-down form as the load increased and as resource utilization grew higher. This, I found, was important so that the application could scale up in case of increased traffic demands while at the same time avoiding costs that are avoidable because they may come with numerous rounds of optimization and testing.

4.4. Solutions to Challenges

- **Proactive Scaling Policies:** With respect to the scaling delays, one effective solution was in the likelihood of the equating of proactive scaling policies. With historical data on traffic patterns, the system was designed to add resources when it anticipated increased traffic intensity. For instance, it was initially configured that the Auto Scaling group should initiate the addition of EC2 instances 10 minutes before the peak hours. Such a proactive measure guaranteed that supplementary instances were running already by the time the traffic demands increased, thereby eliminating instance boot-up times as well as load balancer registration time. Pre-emptive scaling also enabled smooth performance throughout sudden user activity surges that normal scaling measures would take to execute, thus avoiding the performance decline that sticky scaling creates. With cyclic scaling, the system would be able to dynamically prepare for high-traffic periods, as demonstrated in the case where the application was able to respond effectively to higher loads.
- Real-Time Metrics: One of the critical solutions to scale delays and the highly convoluted configuration process was using AWS CloudWatch, which allowed real-time metric metrics. Soon, the system could scale automatically and by controller, constantly tracking currently important parameters, like CPU load, amount of memory, traffic intensity and so on; if the threshold of the limit was exceeded, the system immediately would begin increasing the working nodes. This offered, for the first time, a wall-to-wall, real-time perception of the health/ performance of the application and the ability to scale precisely. Furthermore, available AWS services, namely predictive scaling features, were considered in order to increase the system's capacity for traffic changes. Scalability employs the use of big data and predictive analytics to predict future consumption and then scales based on the expected load for operation within minutes. This significantly reduced the time-consuming traditional reactive scaling strategies and let the system be more prepared to respond to future fluctuating traffic loads. These individual features of real-time monitoring and predictive



scaling made the application performance less sensitive to sudden traffic surges, and scaling became simultaneous, fast and less sensitive to delay or configuration issues.

5. Conclusion

5.1. Summary

This research aims to identify the features, advantages, and importance of using AWS Auto Scaling in current approaches to software development, which will exhibit its critical impact on enhancing the speed, availability, and profitability of applications. The fundamental idea of AWS Auto Scaling is that applications need not be modified when load characteristics change; instead, the resources available can be automatically scaled up or down according to need. These improvements demonstrate the value that Auto Scaling adds by providing, for example, a 73% decrease in response time, a 7.8% increase in the availability of certain resources, and 26.6% overall cost savings when applied to varying traffic situations. This capability to elastically provision resources pays dividends in making sure that applications are ready to deal with both high and low traffic, all in the most efficient use of cloud resources. As illustrated, Auto Scaling not only helps save operational costs but also has a positive impact on the user side since it brings better system response and performs less downtime than having it manually done, which makes it an important tool for organizations that are relying on cloud infrastructures. Additionally, through real-time monitoring and proactive scaling policies mentioned in the study, applications can now be ready to scale when there is a huge surge in traffic as opposed to the usual lagging times in scaling actions. This approach means that resources are indulged before there is a demand for them, which makes performance during high demand very efficient. The continuous monitoring offered through AWS CloudWatch allows the application to quickly address performance issues that slow it down, allowing for immediate action to be taken towards maintaining a highperforming system. The fact that customers are able to scale decisions based on real-time data makes the argument for AWS Auto Scaling for wide-ranging types of applications stronger.

5.2. Future Work

As such, AWS Auto Scaling was successfully applied in this study, but a focus on additional different and even more sophisticated methods for achieving further fine-tuning of the AWS scaling process in future research may be of interest. One of the areas that may show positive signs of change is scaling up with the help of machine learning (ML). Machine learning could generate precise predictions with historical data and traffic to help Auto Scaling not only for current status but for future events as well. The ability to predict scaling activity might prevent unjustified scaling operations, reduce time, and make the scaling process more effective at a lower cost. More specifically, the present study and the resulting set of conclusions might be extended to investigate serverless architectures in the future. Various technologies like AWS Lambda mean serverless architectures that scale applications and compute instances as per the application traffic without worrying about instance management. Extending Endpoints powered by Auto Scaling to serverless architectures would remove the EC2 instance provisioning step altogether and simplify scaling work even more, saving costs and complexity in the process. Future work could explore how Auto Scaling might augment serverless architectures, especially for scale-out to and from other EC2 instances as well as for scale-in from the serverless pile. Finally, the future of AWS Auto Scaling services is integrated with other advanced technologies, such as machine learning and Serverless Applications, to create even more advanced scalable, reliable, and costoptimized application models.



E-ISSN: 2229-7677 • Website: www.ijsat.org • Email: editor@ijsat.org

References

- 1. Wang, X., Xu, C., Wang, K., Yan, F., & Zhao, D. (2020). Memory scaling of cloud-based Big Data systems: A hybrid approach. IEEE Transactions on Big Data, 8(5), 1259-1272.
- 2. Ramachandran, M., & Mahmood, Z. (Eds.). (2020). Software engineering in the era of cloud computing (No. 172538). Berlin/Heidelberg, Germany: Springer.
- 3. Sousa, T. B., Aguiar, A., Ferreira, H. S., & Correia, F. F. (2016, November). Engineering software for the cloud: patterns and sequences. In Proceedings of the 11th Latin-American Conference on Pattern Languages of Programming (pp. 1-8).
- 4. de Sousa, T. B. P. (2020). Engineering Software for the Cloud: A Pattern Language.
- Venugopal, M. V. L. N., & Reddy, C. R. K. (2021). Serverless through cloud native architecture. Int. J. Eng. Res. Technol, 10, 484-496.
- R. Buyya, C. S. Yeo, S. Venugopal, J. Broberg, and I. Brandic, "Cloud Computing and Emerging It Platforms: Vision, Hype, and Reality for Delivering Computing as the 5th Utility," Future Generation Computer Systems, Vol. 25, No. 6, pp. 599-616, (2019).
- 7. Atkinson, C., & Draheim, D. (2013). Cloud-aided software engineering: evolving viable software systems through a web of views. Software engineering frameworks for the cloud computing paradigm, 255-281.
- 8. Kumar, J., & Singh, A. K. (2016, December). Dynamic resource scaling in cloud using neural network and black hole algorithm. In 2016 Fifth International Conference on Eco-friendly Computing and Communication Systems (ICECCS) (pp. 63-67). IEEE.
- Liao, W. H., Kuai, S. C., & Leau, Y. R. (2015, December). Auto-scaling strategy for Amazon web services in cloud computing. In 2015 IEEE International Conference on Smart City/SocialCom/SustainCom (SmartCity) (pp. 1059-1064). IEEE.
- 10. Klein, J., & Van Vliet, H. (2013, June). A systematic review of system-of-systems architecture research. In Proceedings of the 9th international ACM Sigsoft conference on Quality of software architectures (pp. 13-22).
- 11. Jannapureddy, R., Vien, Q. T., Shah, P., & Trestian, R. (2019). An auto-scaling framework for analyzing big data in the cloud environment. Applied Sciences, 9(7), 1417.
- George Fernandez, I., & Arokia Renjith, J. (2021). A Novel Approach on Auto-Scaling for Resource Scheduling Using AWS. In International Virtual Conference on Industry 4.0: Select Proceedings of IVCI4. 0 2020 (pp. 99-109). Springer Singapore.
- 13. da Silva, V. G., Kirikova, M., & Alksnis, G. (2018). Containers for virtualization: An overview. Applied Computer Systems, 23(1), 21-27.
- Gomez-Sanchez, P., Encinas, D., Panadero, J., Bezerra, A., Mendez, S., Naiouf, M., ... & Luque, E. (2016). Using AWS EC2 as Test-Bed infrastructure in theI/O system configuration for HPC applications. Journal of Computer Science and Technology, 16(02), 65-75.
- 15. Reddy, N. A., & Naidu, W. T. Aws Load Balancing Used In Deployments.
- Liao, W.-H., Kuai, S.-C., &Leau, Y.-R. (2019). Auto-scaling Strategy for Amazon Web Services in Cloud Computing. 2019 IEEE International Conference on Smart City/SocialCom/SustainCom (SmartCity). doi:10.1109/smartcity.2019.209.
- 17. Routavaara, I. (2020). Security monitoring in AWS public cloud.



- Rahman, M., Iqbal, S., & Gao, J. (2014, April). Load balancer as a service in cloud computing. In 2014 IEEE 8th international symposium on service oriented system engineering (pp. 204-211). IEEE.
- 19. Harman, M., Lakhotia, K., Singer, J., White, D. R., & Yoo, S. (2013). Cloud engineering is search based software engineering too. Journal of Systems and Software, 86(9), 2225-2241.
- 20. Sarkar, A., & Shah, A. (2018). Learning AWS: Design, build, and deploy responsive applications using AWS Cloud components. Packt Publishing Ltd.