



Placement Empowerment Program

Cloud Computing and DevOps Centre

Implement Auto-scaling in the CloudSet up an autoscaling group for your cloud VMs to handle variable workloads.

Name: Akshaya S

Department: CSE

Introduction

As modern applications face varying workloads, ensuring optimal performance and availability is critical. Auto Scaling, a feature provided by cloud platforms like AWS, dynamically adjusts computing resources in response to demand changes. This Proof of Concept (PoC) demonstrates how to set up an Auto Scaling Group (ASG) for virtual machines (VMs) to handle fluctuating workloads effectively. It explores defining launch configurations, setting scaling policies, and testing automatic scaling based on CPU usage.

Overview

This PoC focuses on implementing a scalable architecture using AWS Auto Scaling Groups. The workflow includes:

- 1. Defining a Launch Template:** Configuring virtual machines (VMs) with required specifications like instance type, AMI, key pairs, and security groups.
- 2. Creating an Auto Scaling Group:** Setting initial group size and linking it to the launch template to manage instances dynamically.
- 3. Configuring Scaling Policies:** Setting up metrics like CPU utilization to trigger scaling actions (e.g., scaling up during high CPU usage).
- 4. Testing Auto Scaling:** Simulating high CPU load to verify that the ASG launches additional instances to handle demand.

This PoC will demonstrate the reliability, flexibility, and costefficiency of dynamic scaling in a cloud environment.

Objective

The primary objective of this PoC is to:

1. Implement an **Auto Scaling Group (ASG)** to manage workloads effectively.
2. Define and configure a **Launch Template** for virtual machines.
3. Set up and test **scaling policies** based on predefined metrics, such as CPU utilization.
4. Validate the scaling process by simulating real-world scenarios (e.g., high CPU usage).

By completing this PoC, the goal is to gain hands-on experience with Auto Scaling and to understand its importance in ensuring application availability and cost management.

Importance

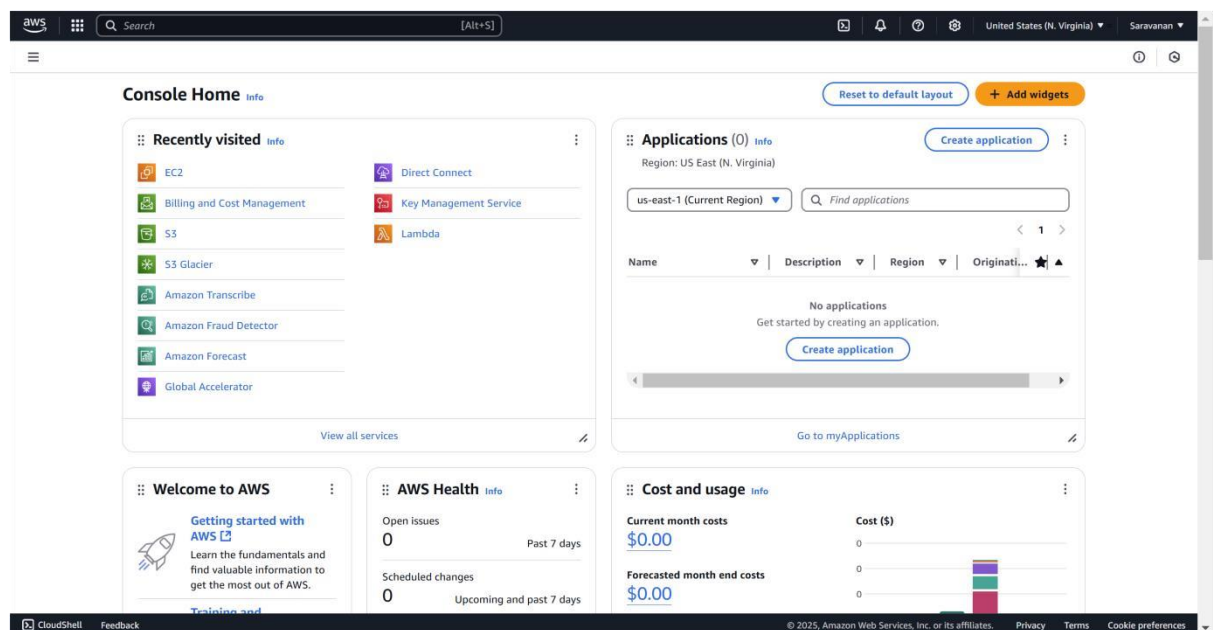
1. **Improved Application Availability:** Auto Scaling ensures that applications remain available even during traffic spikes by automatically adding more VMs to meet demand.
2. **Cost Optimization:** It dynamically reduces the number of VMs during low traffic periods, minimizing unnecessary costs.
3. **Efficient Resource Utilization:** By scaling resources based on actual demand, Auto Scaling prevents over-provisioning and underutilization.
4. **Resilience to Failures:** Auto Scaling can replace unhealthy instances automatically, ensuring consistent application performance.

5. Real-World Relevance: The ability to manage variable workloads is a critical skill in cloud computing and aligns with industry practices.

Step-by-Step Overview

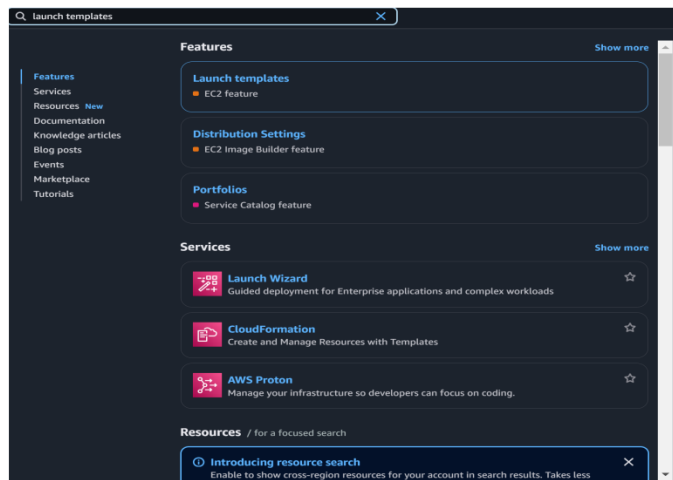
Step 1:

1. Go to [AWS Management Console](#).
2. Enter your username and password to log in.



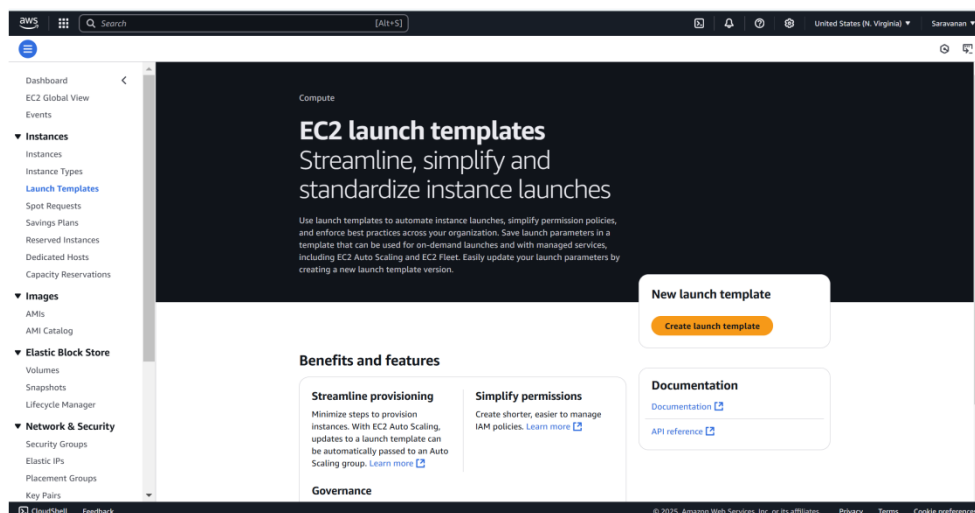
Step 2:

Search for Launch Templates.



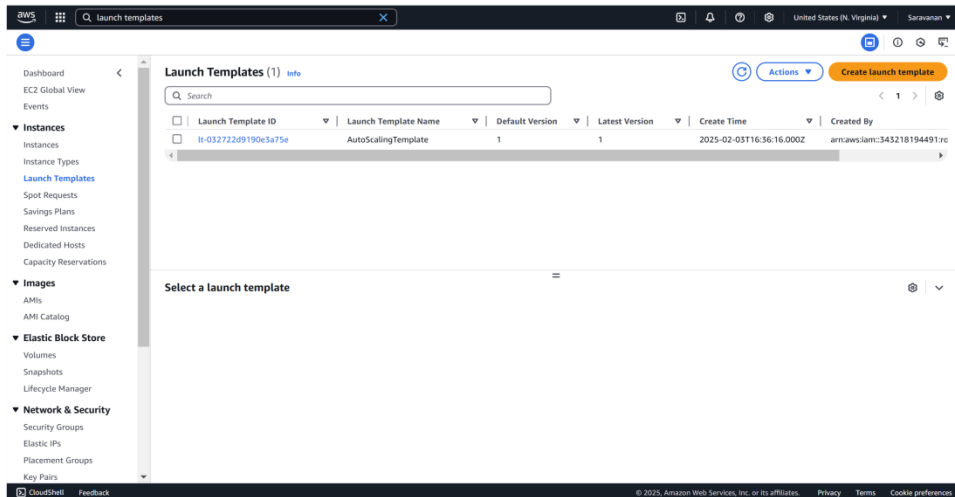
Step 3:

Click on the Create launch template.



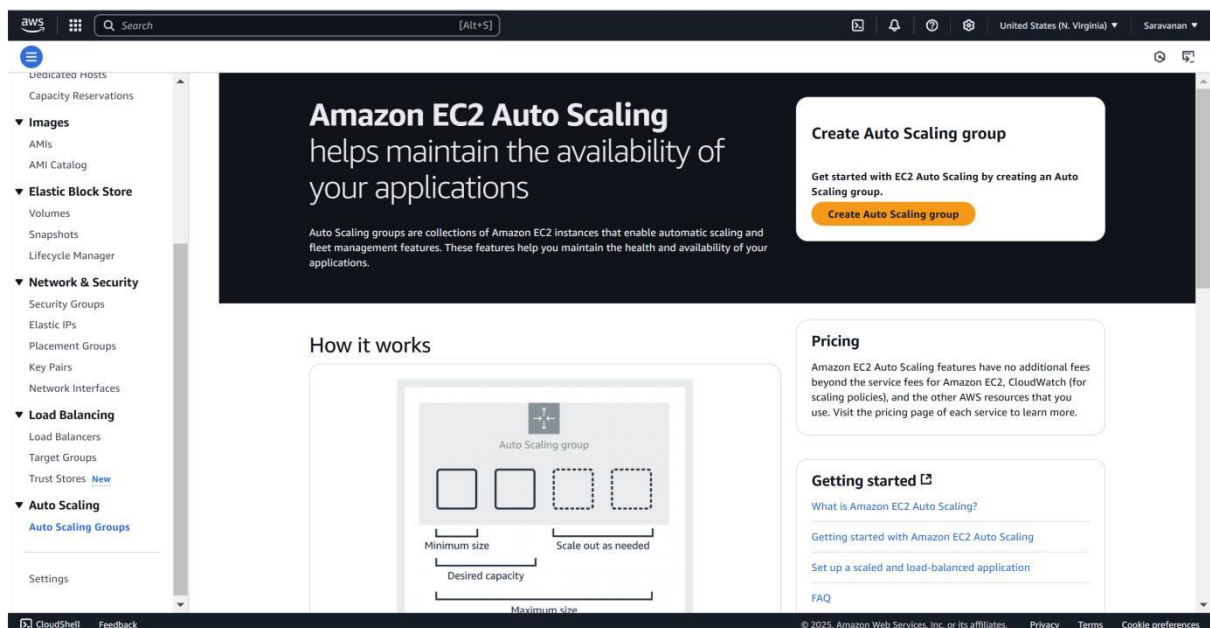
Step 4:

Create a **Launch Template** named **AutoScalingTemplate** using an **Amazon Machine Image (AMI)** like Amazon Linux 2 or any default image, and choose an **instance type** such as **t2.micro** for free-tier eligibility. Select an **existing key pair** (or create a new one) to enable SSH access, and configure a **security group** that allows HTTP (port 80) and SSH (port 22). Once all details are filled out, click **Create launch template** to complete the setup.



Step 5:

Go to the **EC2 Dashboard**. On the left sidebar, click on **Auto Scaling Groups**. Click on **Create an Auto Scaling group**.



Step 6:

Auto Scaling group name: Give it a name (e.g., MyAutoScalingGroup).

Launch Template: Select the launch template you created earlier (AutoScalingTemplate).

The screenshot shows the 'Choose launch template' step in the AWS console. The left sidebar lists steps from 'Choose launch template' to 'Review'. The main content area has a 'Name' section with a text input 'MyAutoScalingGroup'. Below is a 'Launch template' section with a dropdown menu showing 'AutoScalingTemplate' and a 'Version' dropdown showing 'Default (1)'. A 'Description' field is empty. At the bottom, there are links for 'Launch template', 'Instance type', 'AMI ID', 'Security groups', and 'Resource Spot Instances'.

Step 7:

VPC and Subnets: Choose your VPC (it's fine to use the default one). Select at least two subnets in different Availability Zones (this ensures high availability).

The screenshot shows the 'Choose instance launch options' step in the AWS console. The left sidebar lists steps from 'Choose instance launch options' to 'Review'. The main content area has an 'Instance type requirements' section with a table showing 'Launch template' as 'AutoScalingTemplate', 'Version' as 'Default', and 'Description' as '-'. Below this is a 'Network' section with a 'VPC' dropdown showing 'vpc-0f36f0944c12862e5' and 'Availability Zones and subnets' dropdown showing 'us-east-1a | subnet-0f01daeb9f48d5fc4'. There are links to 'Override launch template', 'Create a VPC', and 'Create a subnet'.

Step 8:

For this PoC leave the next settings as default and click next .

 CloudShell [Feedback](#) © 2025 Amazon Web Services, Inc. or its affiliates [Privacy](#) [Terms](#) [Cookie preferences](#)

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Step 9:

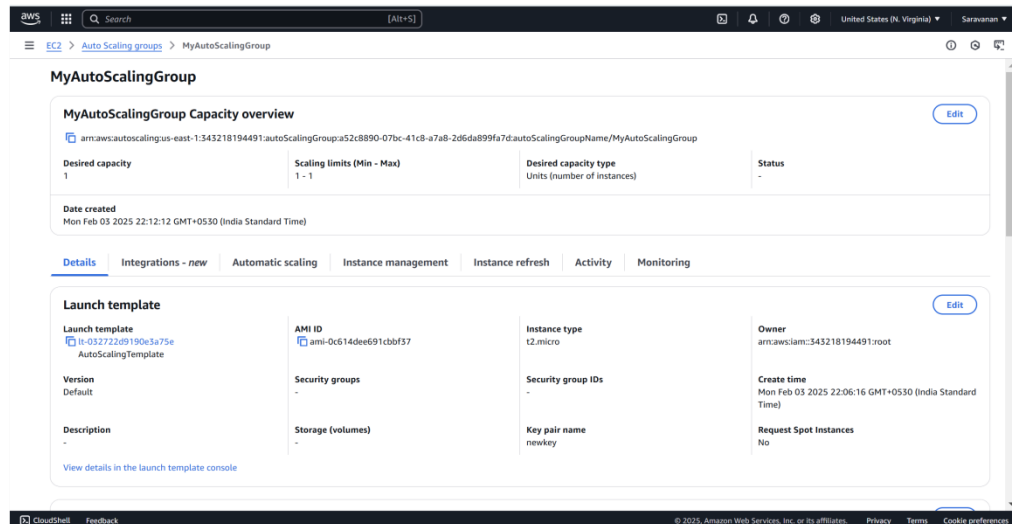
Review all the settings you've configured. Once satisfied, click **Create Auto Scaling Group**.

The screenshot shows the 'Review' step of the 'Create Auto Scaling group' wizard in the AWS Management Console. On the left, a progress bar lists seven steps: 1. Choose launch template, 2. Choose instance launch options, 3 - 6 (optional), and 7. Review. Step 7 is currently selected. The main content area is divided into two sections: 'Step 1: Choose launch template' and 'Step 2: Choose instance launch options', each with an 'Edit' button. The 'Group details' section shows the 'Auto Scaling group name' as 'MyAutoScalingGroup'. The 'Launch template' section shows 'Launch template' as 'AutoScalingTemplate' and 'Version' as 'Default'. The 'Network' section shows 'VPC' as 'vpc-0f36f0944c12862e5'. The 'Availability Zones and subnets' table lists 'us-east-1a' with subnet 'subnet-0f01d5eb9f48d5fc4' and CIDR range '172.31.80.0/20'. The 'Availability Zone distribution' is set to 'Balanced best effort'.

Availability Zone	Subnet	Subnet CIDR range
us-east-1a	subnet-0f01d5eb9f48d5fc4	172.31.80.0/20

The screenshot shows the 'Auto Scaling groups' page in the AWS Management Console. At the top, there are tabs for 'Launch configurations', 'Launch templates', 'Actions', and a prominent 'Create Auto Scaling group' button. Below the tabs is a search bar and a table of Auto Scaling groups. The table has columns for Name, Launch template/configuration, Instances, Status, Desired capacity, Min, Max, and Availability Zones. One group, 'MyAutoScalingGroup', is listed with launch template 'AutoScalingTemplate', version 'Default', 0 instances, and a status of 'Updating capacity...'. The bottom of the page shows '0 Auto Scaling groups selected'.

Name	Launch template/configuration	Instances	Status	Desired capacity	Min	Max	Availability Zones
MyAutoScalingGroup	AutoScalingTemplate Version Default	0	Updating capacity...	1	1	1	us-east-1a



Step 10:

Testing Auto Scaling :

Important Note

Do Not Perform This Test If You Want to Avoid Costs:

1. Launching and running additional EC2 instances will incur charges beyond the AWS Free Tier.
2. Simulating high CPU usage and triggering scaling may increase costs temporarily due to additional resource allocation.

1. Simulate High CPU Usage on an EC2 Instance

Connect to one of your EC2 instances in the Auto Scaling Group using SSH.

Run a command to create artificial CPU load. For example:

```
sudo yum install -y stress
```

```
stress --cpu 2 --timeout 300
```

This command will utilize 2 CPU cores for 5 minutes, simulating high CPU usage.

2. Monitor Scaling Activities

Navigate to the **AWS Management Console > EC2 Dashboard > Auto Scaling Groups**.

Select your Auto Scaling Group and go to the **Activity History** tab.

Check if a new instance is being launched based on your scaling policy (e.g., CPU utilization exceeding 50%).

3. Terminate the Stress Test

Once testing is done, stop the CPU load by pressing Ctrl+C in the terminal or by terminating the stress process.

4. Verify Scaling Down

After the CPU usage drops, monitor the Auto Scaling Group again to confirm that unnecessary instances are terminated, returning to the desired capacity.