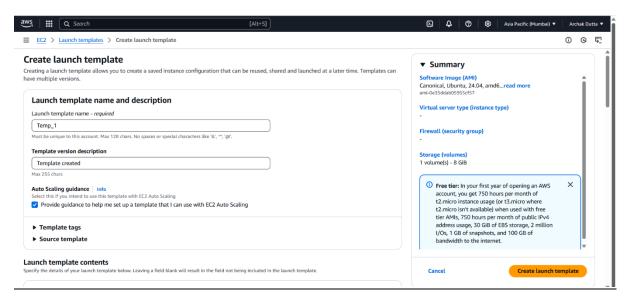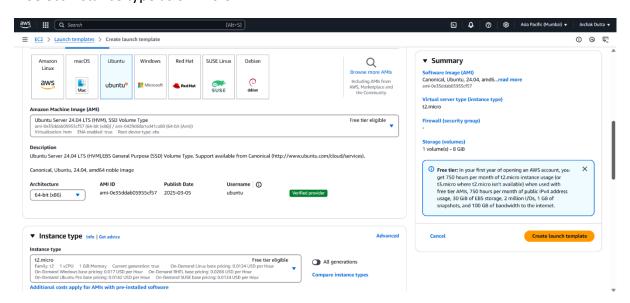*Assignment No: 11*

*Problem Statement: **Build scaling plans in AWS that balance the load on different EC2 instances.***

*Solution:* To build different scaling plans in AWS, the steps are-

1.Sign into AWS console, go to EC2 and select **launch template**.Provide the **template name** along with a **description** and the auto scaling option is to be checked.Then select **ubuntu** as OS in quick start.
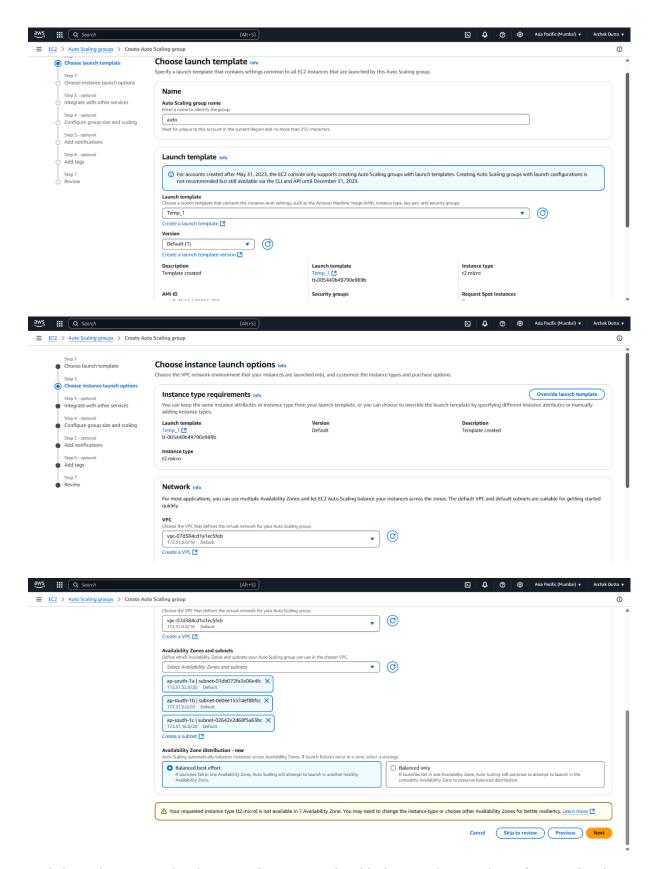


2.Select instance type as **t2.micro**.



3.Select a key pair along with the existing security group from network section.Provide the user data under additional settings.Select create launch template and the template is created.

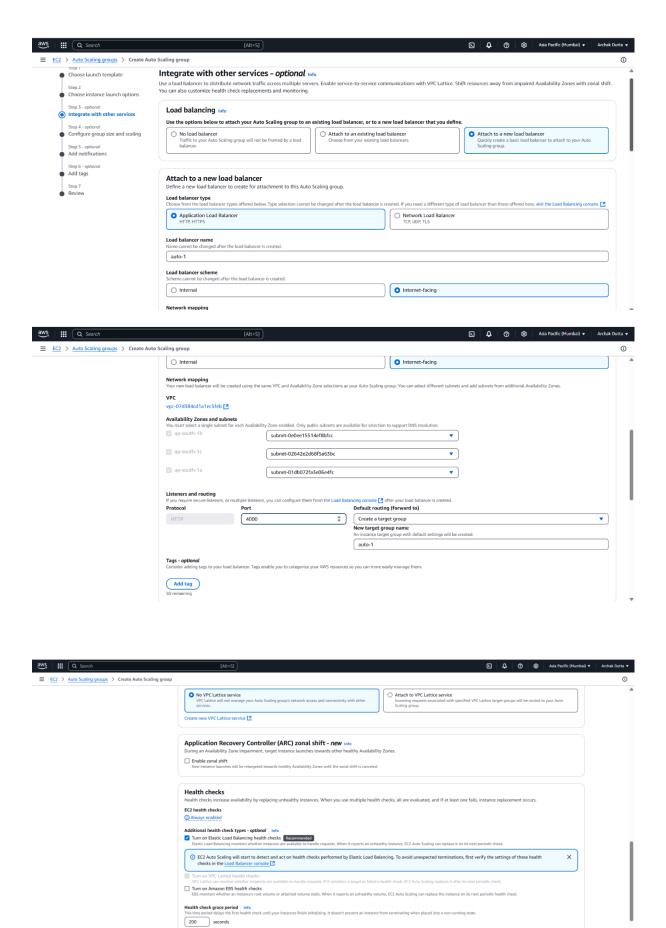4.Select auto scaling group.Provide a **name**,**select a temp**late,**select the zones** under network section and choose **the balanced best effort** option.

EC2 > Auto Scaling groups > Create Auto Scaling group

**Step 1**
Choose launch template

**Step 2**
Choose instance launch options

**Step 3 - optional**
Integrate with other services

**Step 4 - optional**
Configure group size and scaling

**Step 5 - optional**
Add notifications

**Step 6 - optional**
Add tags

**Step 7**
Review

## Choose launch template Info

Specify a launch template that contains settings common to all EC2 instances that are launched by this Auto Scaling group.

### Name

**Auto Scaling group name**
Enter a name to identify the group.

auto

Must be unique to this account in the current Region and no more than 255 characters.

### Launch template Info

ⓘ For accounts created after May 31, 2023, the EC2 console only supports creating Auto Scaling groups with launch templates. Creating Auto Scaling groups with launch configurations is not recommended but still available via the CLI and API until December 31, 2023.

**Launch template**
Choose a launch template that contains the instance-level settings, such as the Amazon Machine Image (AMI), instance type, key pair, and security groups.

Temp_1 ▾

Create a launch template ↗

**Version**

Default (1) ▾

Create a launch template version ↗

| Description | Launch template | Instance type |
|---|---|---|
| Template created | Temp_1 ↗ | t2.micro |
| | lt-005440b49790e989b | |

**AMI ID** | **Security groups** | **Request Spot Instances**

---

## Choose instance launch options Info

Choose the VPC network environment that your instances are launched into, and customize the instance types and purchase options.

### Instance type requirements Info

Override launch template

You can keep the same instance attributes or instance type from your launch template, or you can choose to override the launch template by specifying different instance attributes or manually adding instance types.

**Launch template**
Temp_1 ↗
lt-005440b49790e989b

**Version**
Default

**Description**
Template created

**Instance type**
t2.micro

### Network Info

For most applications, you can use multiple Availability Zones and let EC2 Auto Scaling balance your instances across the zones. The default VPC and default subnets are suitable for getting started quickly.

**VPC**
Choose the VPC that defines the virtual network for your Auto Scaling group.

vpc-07d384cd1a1ec5feb
172.31.0.0/16    Default

Create a VPC ↗

---

Choose the VPC that defines the virtual network for your Auto Scaling group.

vpc-07d384cd1a1ec5feb
172.31.0.0/16    Default

Create a VPC ↗

**Availability Zones and subnets**
Define which Availability Zones and subnets your Auto Scaling group can use in the chosen VPC.

Select Availability Zones and subnets ▾

ap-south-1a | subnet-01db072fa3e06e4fc ✕
172.31.32.0/20    Default

ap-south-1b | subnet-0e0ee15514ef8bfcc ✕
172.31.0.0/20    Default

ap-south-1c | subnet-02642e2d68f5a63bc ✕
172.31.16.0/20    Default

Create a subnet ↗

**Availability Zone distribution - new**
Auto Scaling automatically balances instances across Availability Zones. If launch failures occur in a zone, select a strategy.

🔘 **Balanced best effort**
If launches fail in one Availability Zone, Auto Scaling will attempt to launch in another healthy Availability Zone.

○ **Balanced only**
If launches fail in one Availability Zone, Auto Scaling will continue to attempt to launch in the unhealthy Availability Zone to preserve balanced distribution.

⚠ Your requested instance type (t2.micro) is not available in 1 Availability Zone. You may need to change the instance type or choose other Availability Zones for better resiliency. Learn more ↗
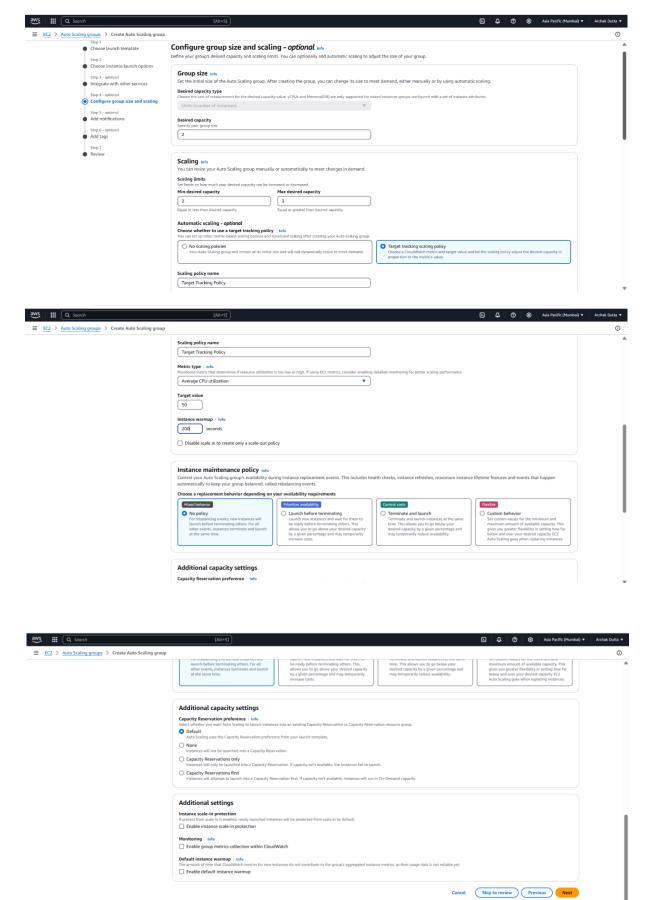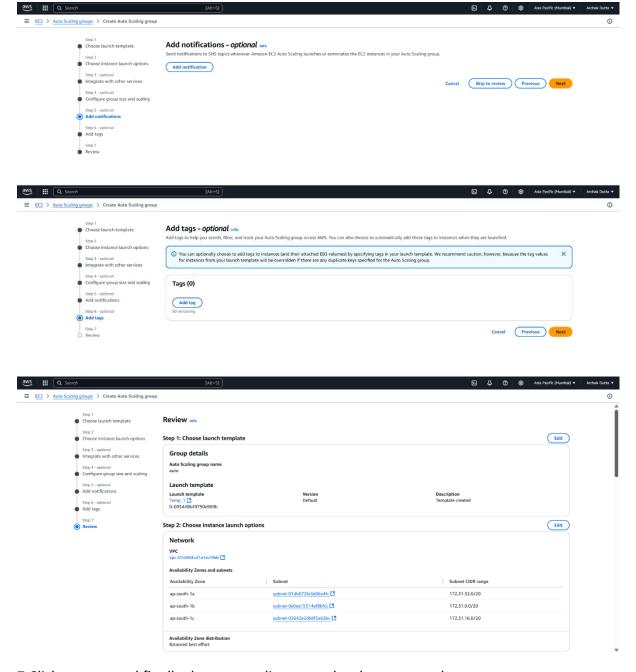
Cancel    Skip to review    Previous    Next

---

5.Click on the next and select **attach to a new load balancer** along with **application load balancer** as **load balancer type** and internet-facing as **the load balancer scheme**.Provide port as **4000** and select create a target group under **listeners and routing**.Finally select **no**

**VPC lattice service** option and enable the additional health check types.Provide **health check grace period as 200 seconds.**
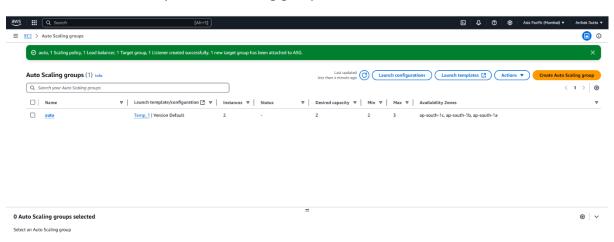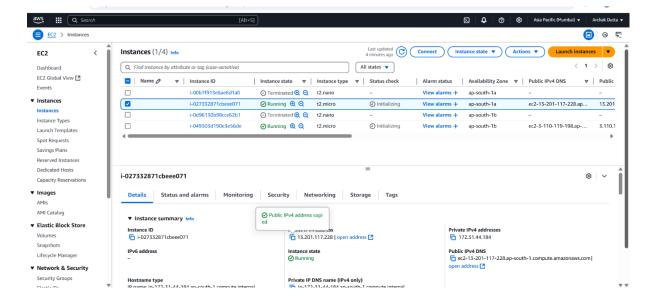
6.In configure group and scaling window, select **desired capacity** as 2.Now under scaling select **min and max desired capacity** as 2 and 3 respectively.Select **target tracking scaling policy** under **automatic scaling**.Now provide **instance warmup as 200 seconds**.
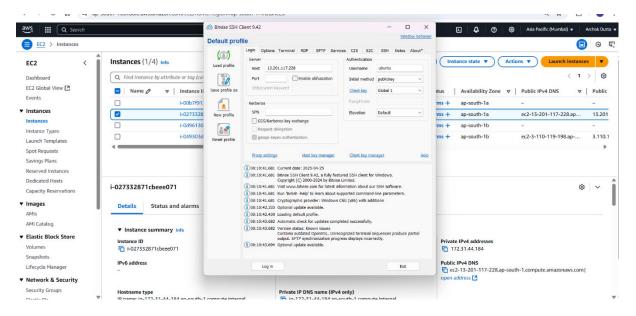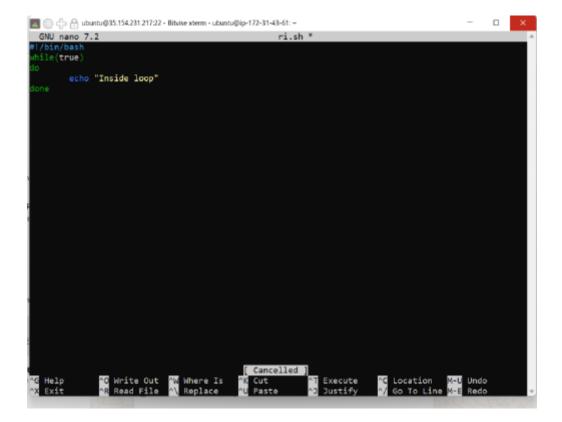
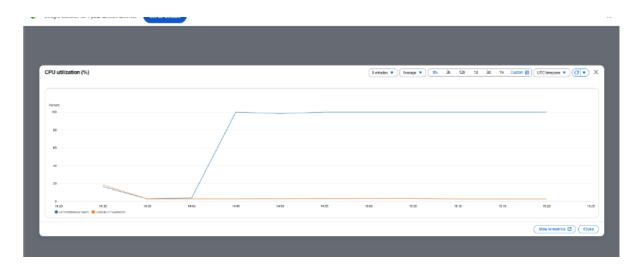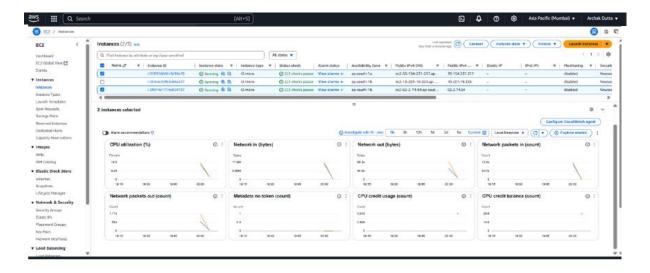7.Click on next and finally the auto scaling group has been created.

8.Open the bitvise ssh client and log in using the given key pair and execute the following commands in the terminal.

9. Now check the CPU utilization on the monitoring section after executing the bash file and the third instance is created as a result of it.

The commands used in this assignment are-

```bash
#!/bin/bash

apt-get update

apt-get install -y nginx

systemctl start nginx

systemctl enable nginx

apt-get install -y git

curl -SL https://deb.nodesource.com/setup_16.x|sudo -E bash -

apt-get install -y nodejs

git clone http://github.com/sudip7407/Repo1.git

cd Repo1

npm install

node index.js
```

//Sudo nano ri.sh

Within this file we will write the following code-

```bash
#!/bin/bash

while(true)

do

    echo "Inside loop"

done
```

//sudo chmod +x ri.sh

//sh ri.sh