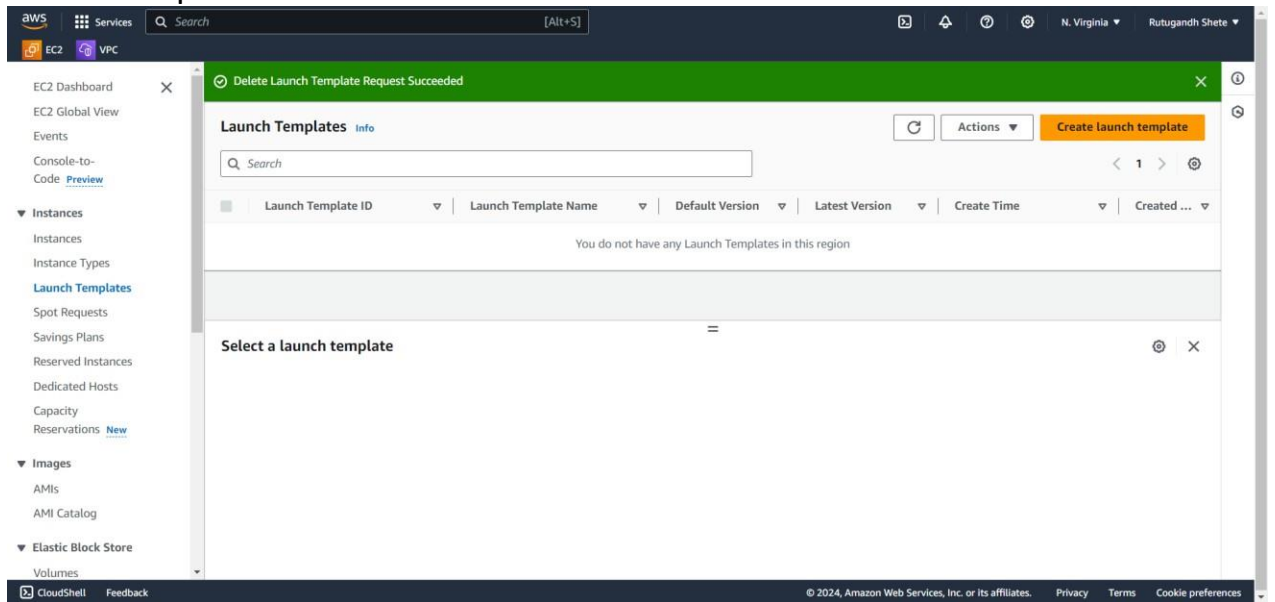| No:- | Content |
|------|---------|
| **1.** | Auto scaling |

Auto Scaling:

Amazon EC2 Auto Scaling helps you ensure that you have the correct number of Amazon EC2 instances available to handle the load for your application. You create collections of EC2 instances, called *Auto Scaling groups*. You can specify the minimum number of instances in each Auto Scaling group, and Amazon EC2 Auto Scaling ensures that your group never goes below this size. You can specify the maximum number of instances in each Auto Scaling group, and Amazon EC2 Auto Scaling ensures that your group never goes above this size. If you specify the desired capacity, either when you create the group or at any time thereafter, Amazon EC2 Auto Scaling ensures that your group has this many instances. If you specify scaling policies, then Amazon EC2 Auto Scaling can launch or terminate instances as demand on your application increases or decreases.

For example, the following Auto Scaling group has a minimum size of four instances, a desired capacity of six instances, and a maximum size of twelve instances. The scaling policies that you define adjust the number of instances, within your minimum and maximum number of instances, based on the criteria that you specify.
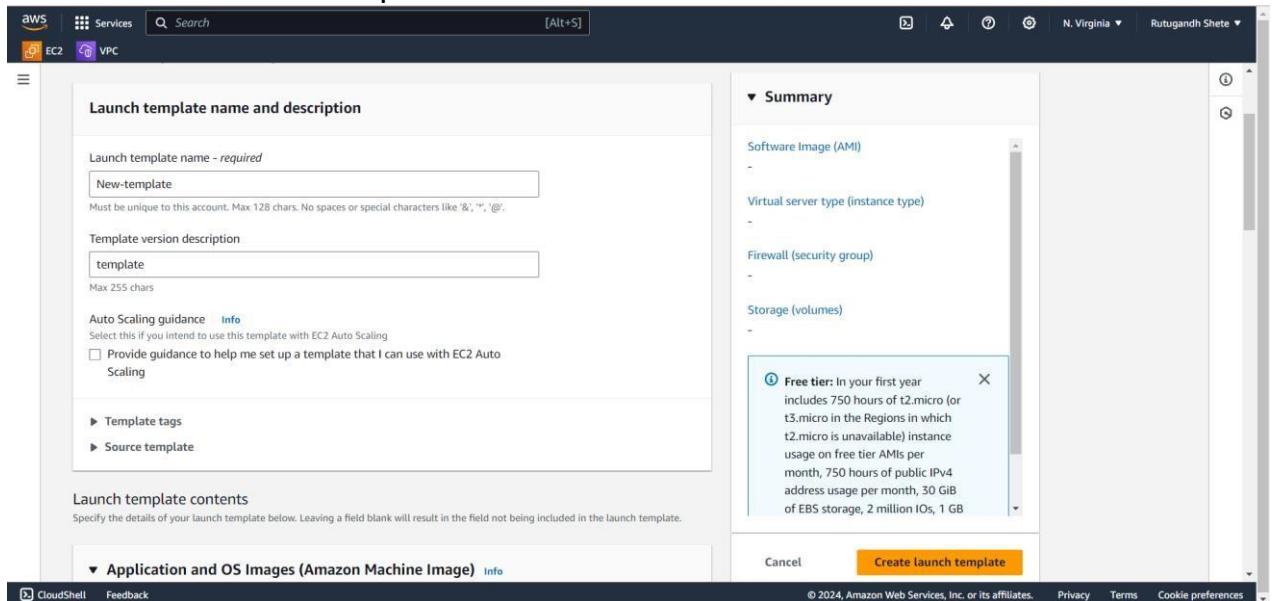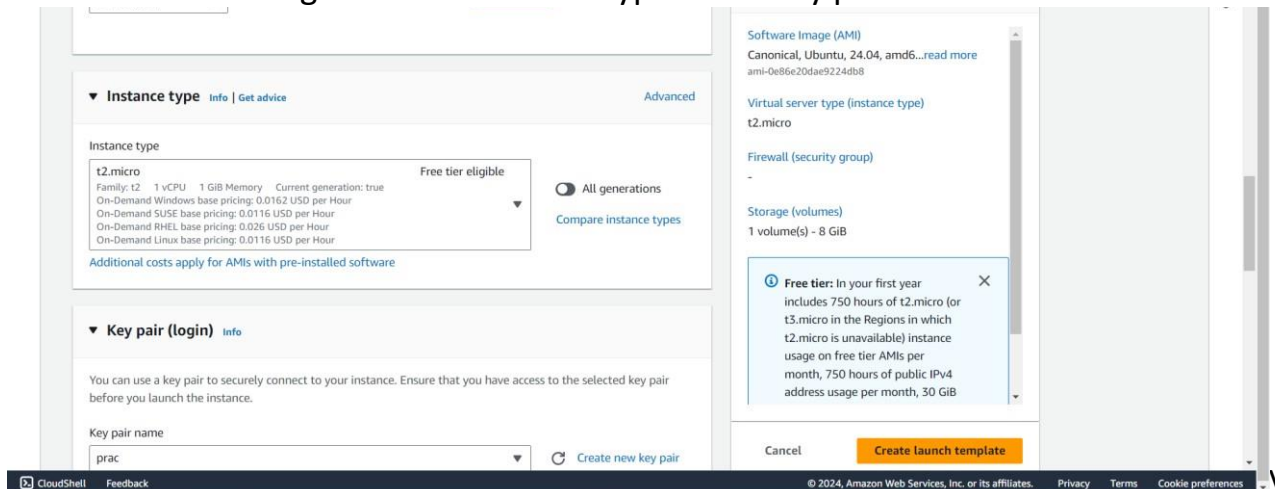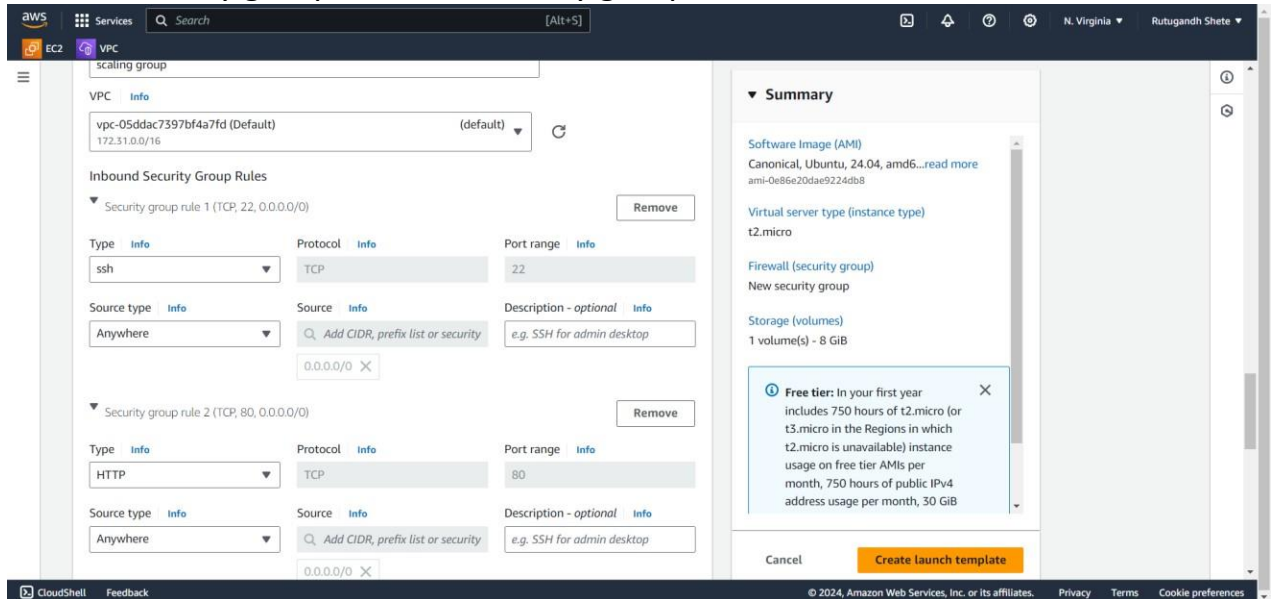
Steps:

- Create template
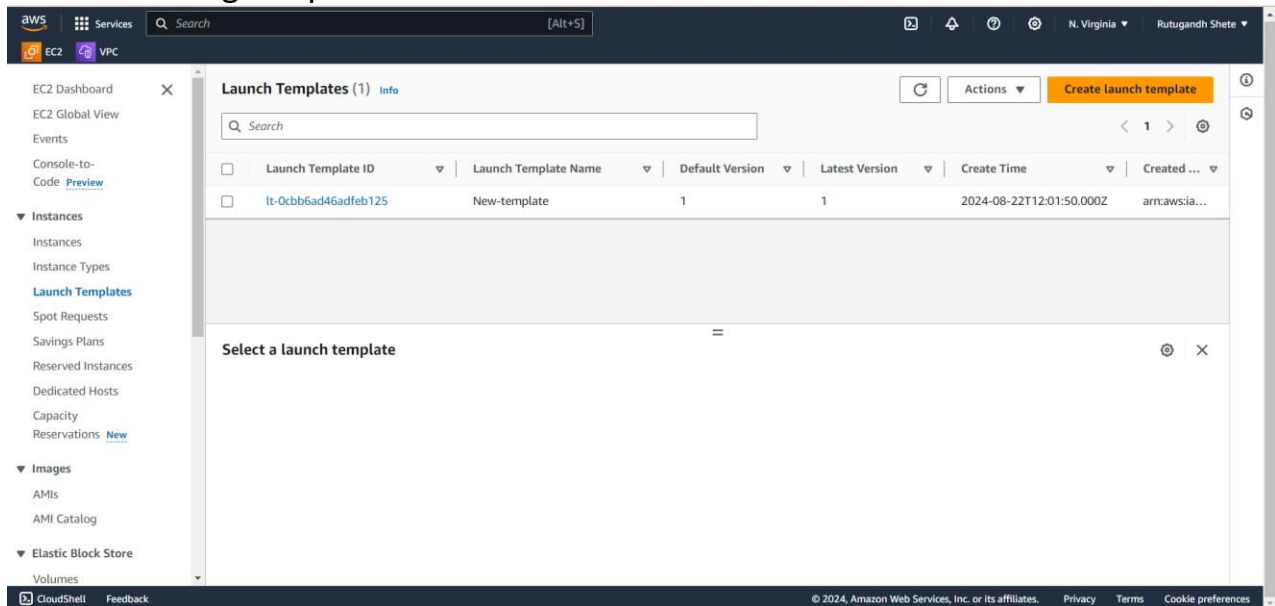


- Provide name and description



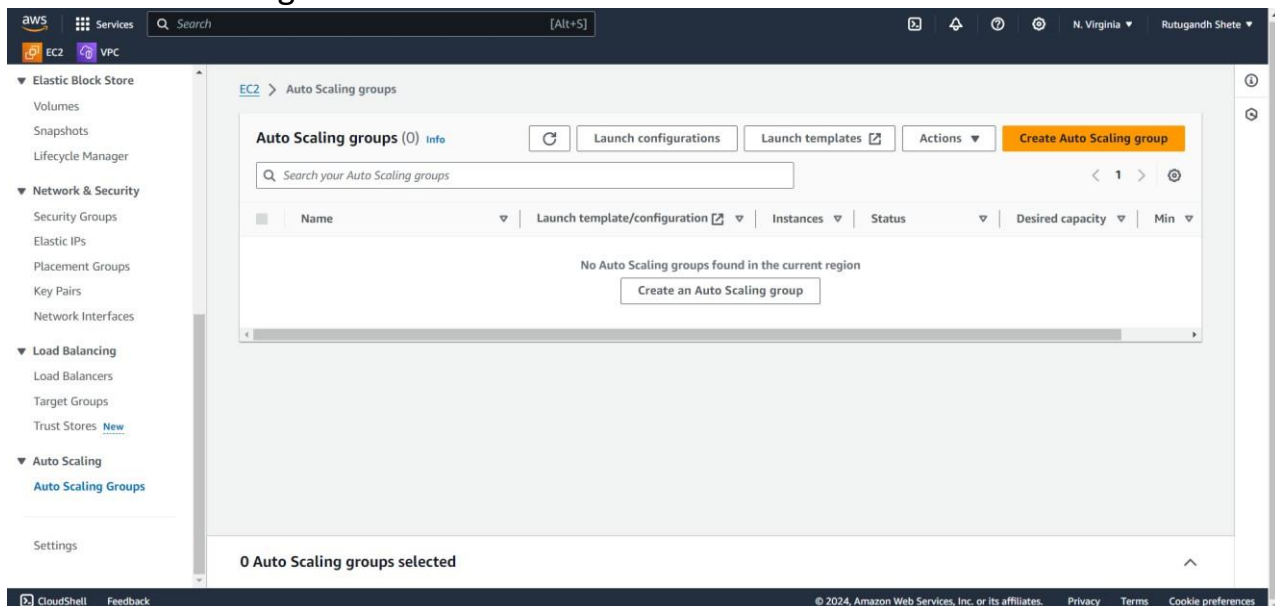- Select machine image→select instance type→add key pair value

- Create security group and add security groups rule.



- After launching template.



- Create auto scaling

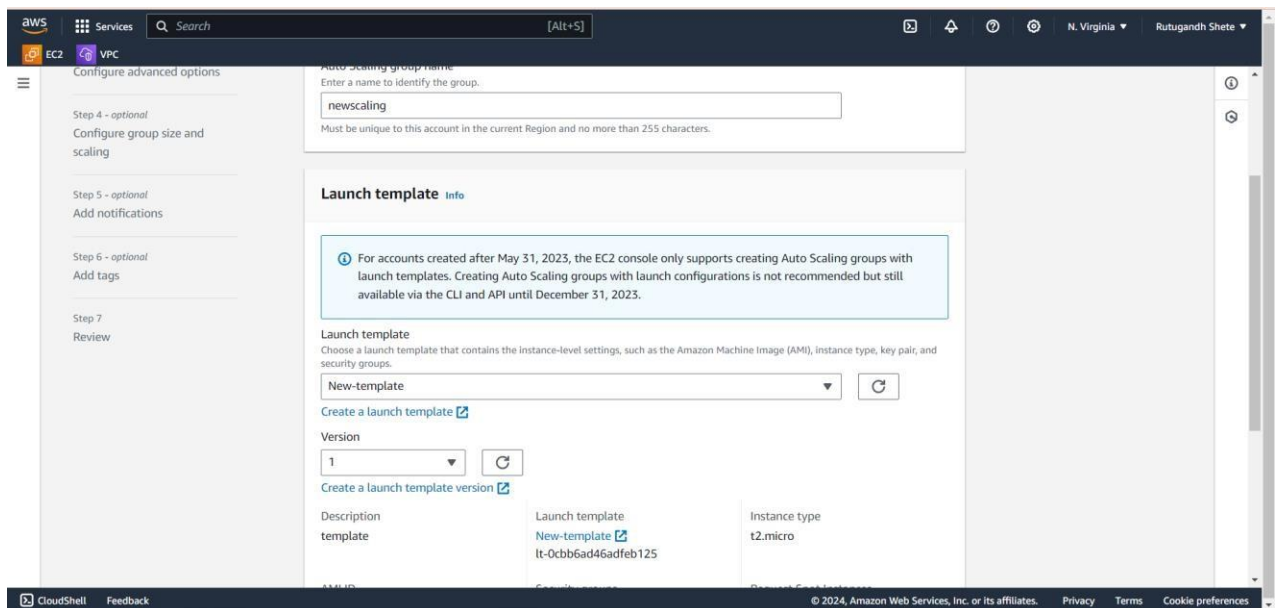- Add name→select template that we have created→select version
  **Default Version:**
  When you create a launch template, AWS automatically designates the first version as the "default version."
  The default version is used by the Auto Scaling group unless you specify a different version.
  **Specifying a Version:**
  When configuring your Auto Scaling group, you can specify which version of the launch template to use.
  If you don't specify a version, the Auto Scaling group will use the default version.



- After selecting availability zones→next

- Select desired capacity → in scaling set min and max limits



- Click on target tracking scaling policy-→add name→add metric type(CPU utilization)→target value →instance warmup(ex- 50)



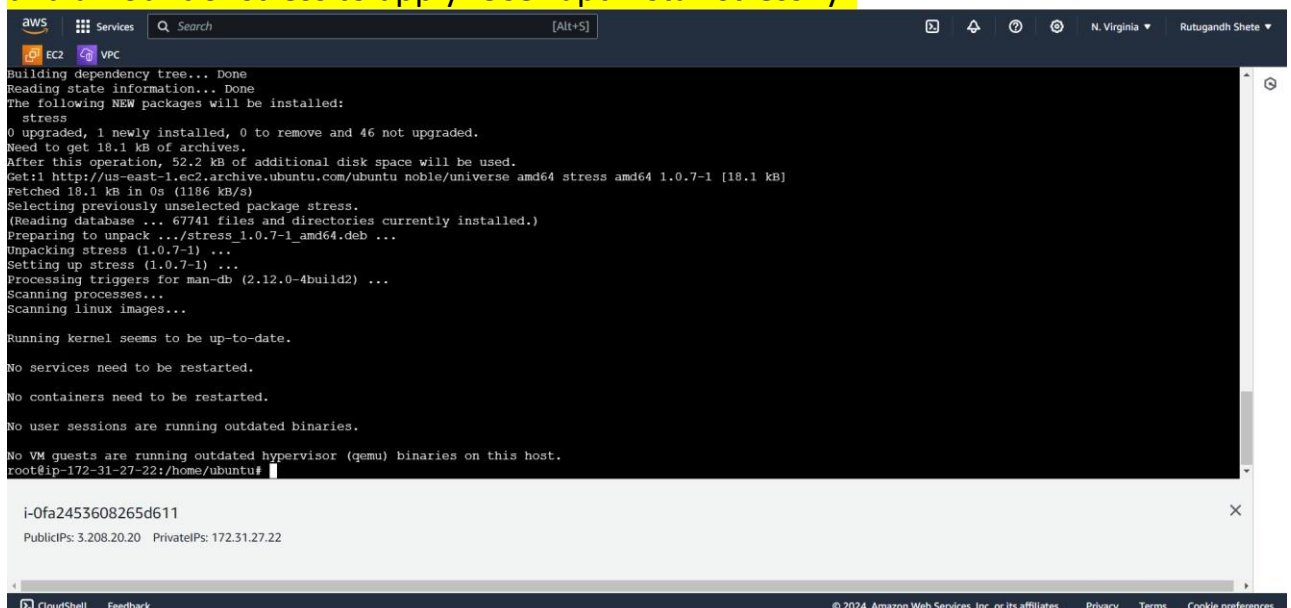- After reviewing details →create auto scaling

- As we can see one instance is created because our min desired capacity is "1".



- Connect that to terminal



- Install stress command. The stress command has several options to specify the type and amount of stress to apply. Use "apt install stress -y"

- Use stress command to see details.



- After adding " stress --cpu 90 --io 4 --vm 2 --vm-bytes 128M --timeout 10m"

--cpu 90: This will start 90 CPU stressor processes. Each process will try to fully utilize one CPU core. This is quite intensive, especially if you have fewer than 90 cores available, as it will overburden the CPU.

io 4: This starts 4 I/O stressor processes. These processes will perform continuous I/O operations (reading and writing data), stressing the system's input/output subsystem.

vm 2: This starts 2 virtual memory stressor processes. These processes will continuously allocate, use, and free memory.

vm-bytes 128M: This sets the amount of memory allocated by each --vm process to 128 MB. Each of the 2 --vm processes will allocate and work with 128 MB of memory.

timeout 10m: This sets the duration of the stress test to 10 minutes. After 10 minutes, the stress command will automatically stop.
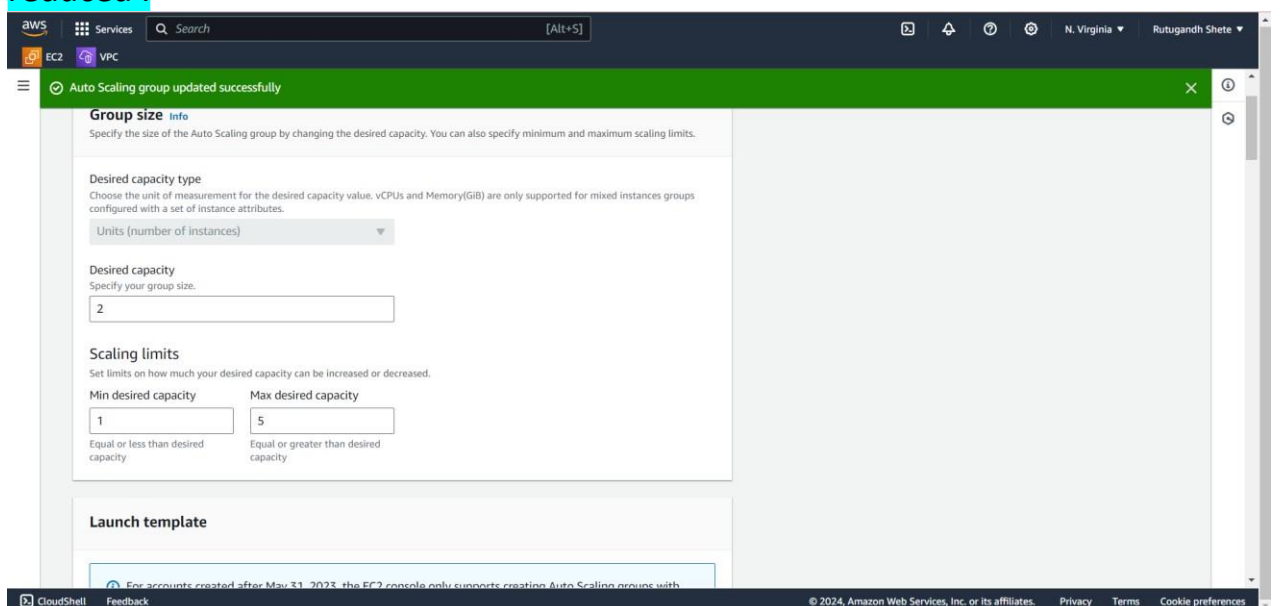
Use "top" to see CPU utilization

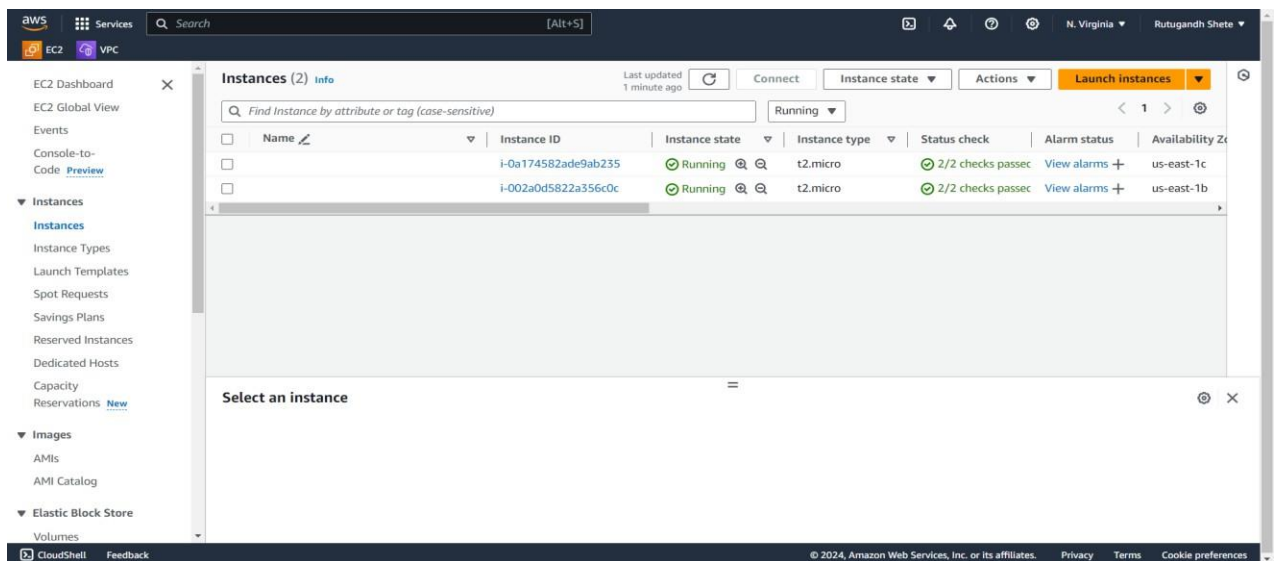- As we can one instance is created after applying stress.



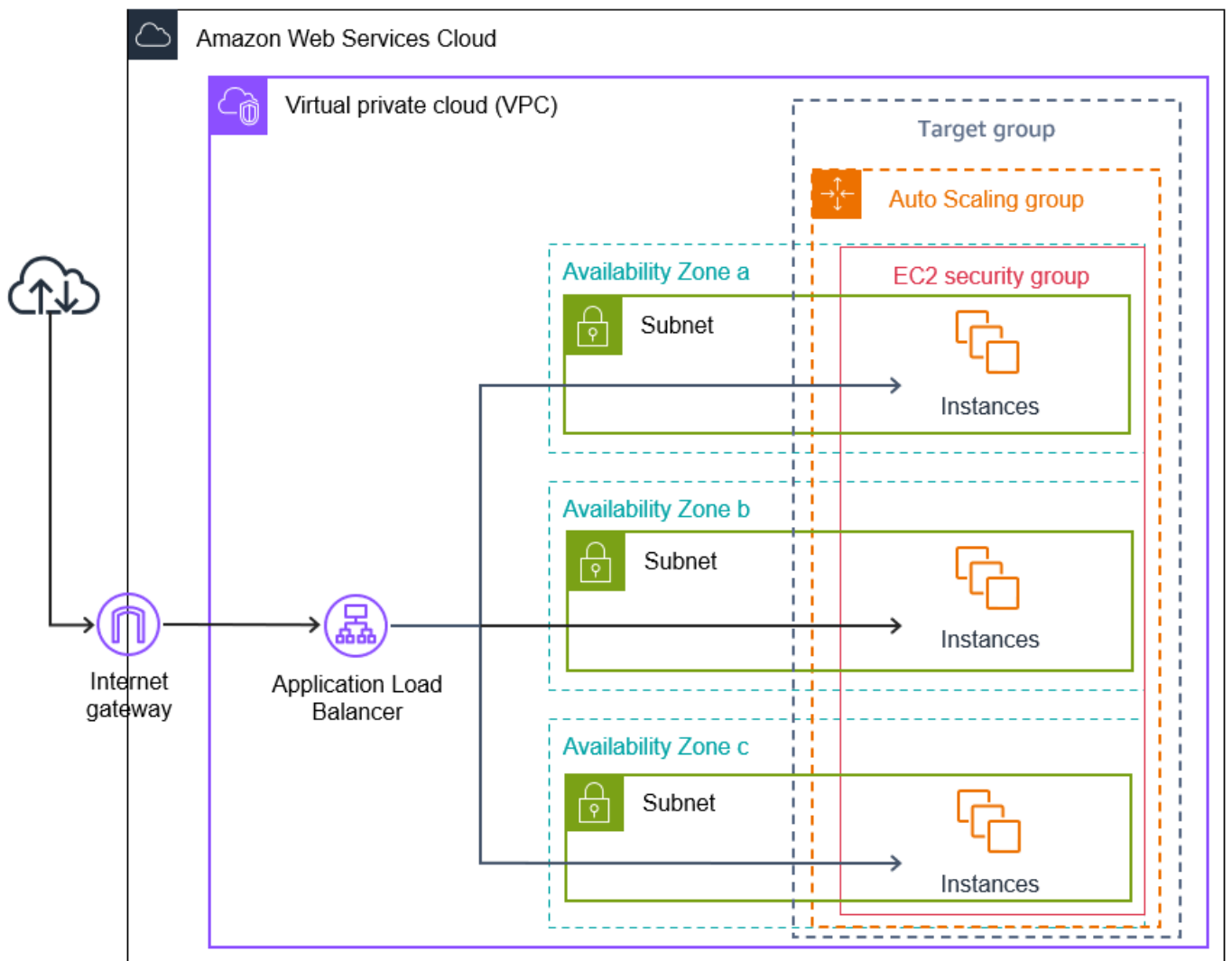- As we can see applying stress to 90% cpu utilization we have achieved "max desired capacity"



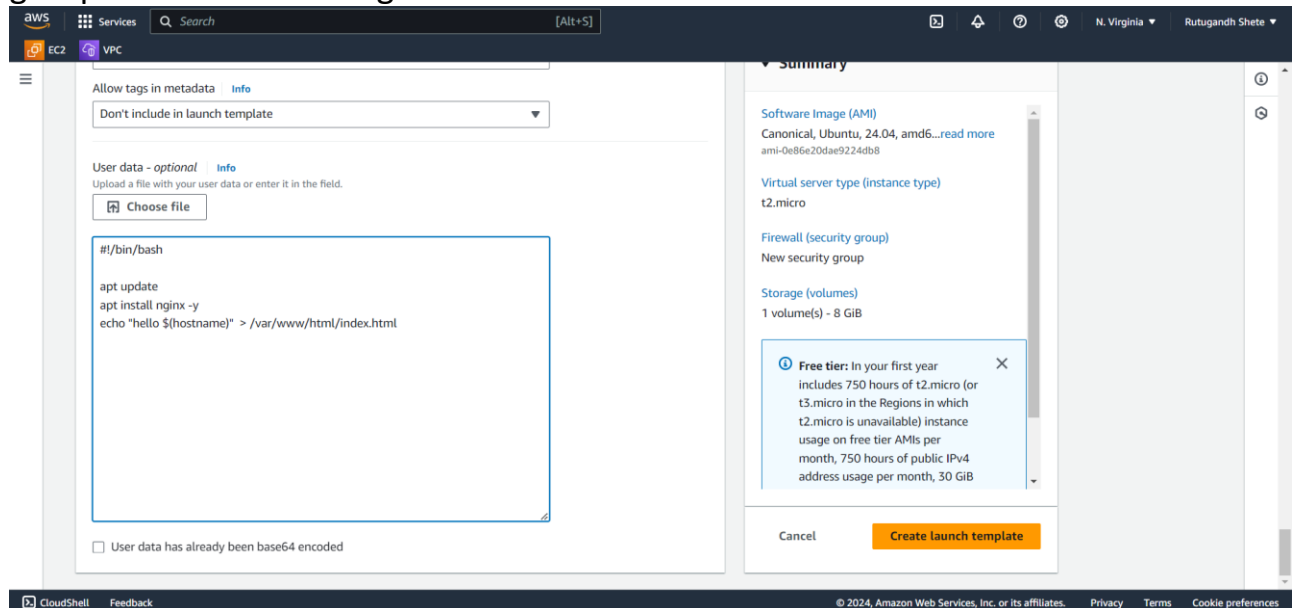- After editing desired capacity we can scale down and number of instances will be reduced .
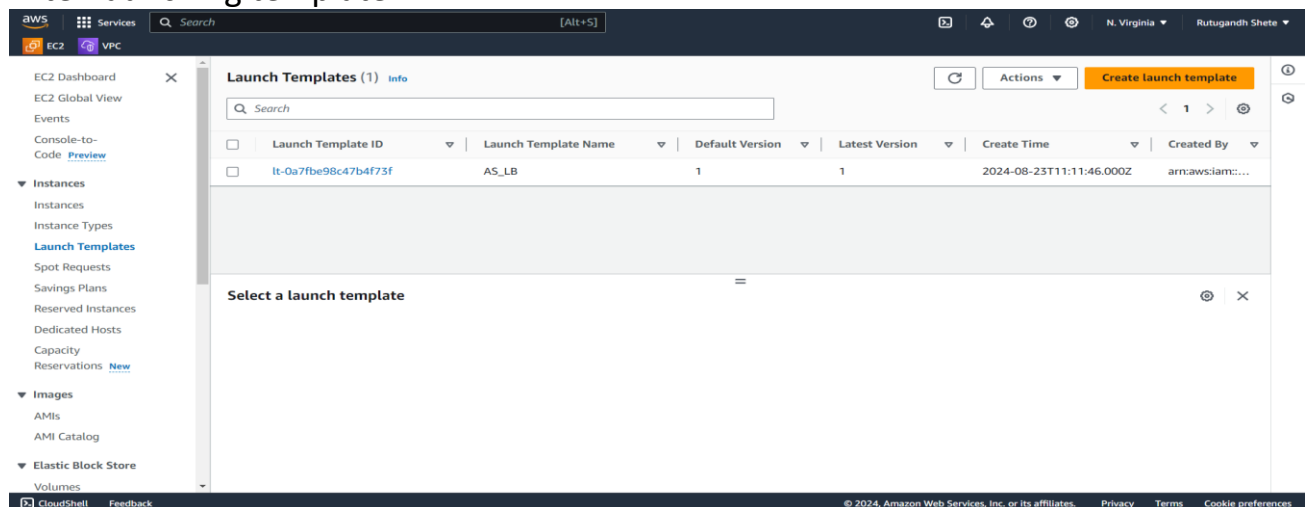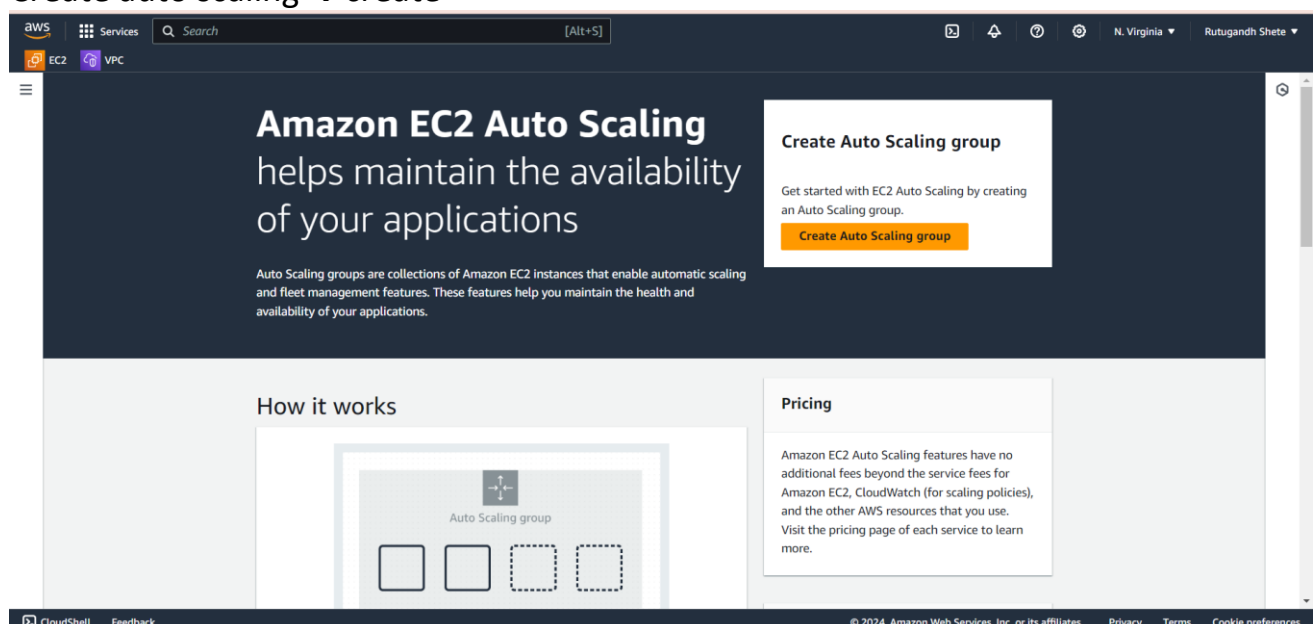
Auto scaling using load balancer:

Steps:

- Create template→Name→AMI→Key pair→instance type→create security group→advanced settings



- After launching template



- Create auto scaling →create

- Provide name →template that we have created →select version



- Select VPC and zones

- Load balancer→select that we have to create new load balancer (if have load balancer then we can choose)→select load balancer type→select Health check grace period (30sec).

  An internal load balancer is only accessible within your VPC. It doesn't have a publicly resolvable DNS name, which means it's not accessible from the internet.

  An internet-facing load balancer has a publicly resolvable DNS name, making it accessible from the internet.



- Create target →name

- Provide desired, min and max desired capacity.



- Add target scaling policy→name→metric type→target value→instance warmup.



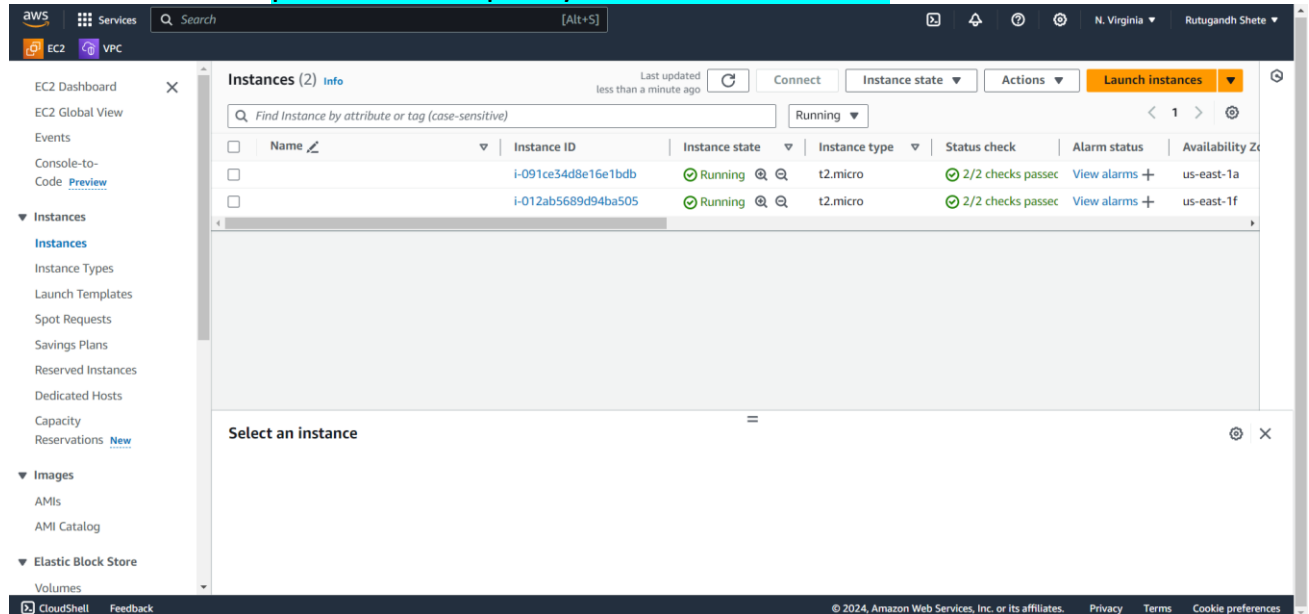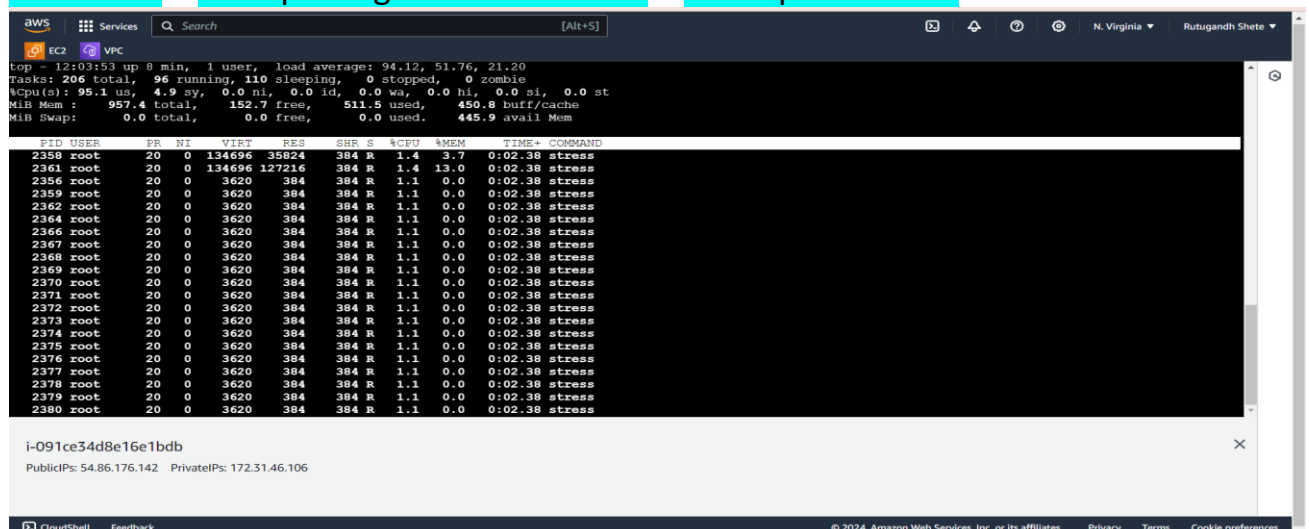- Create auto scaling group

- As we can see that our load balancer is in active state and paste DNS onto browser



- As we can see as per desired capacity we have 2 instances.



- Connect instance to teminal →install stress command(apt install stress -y)→use stress command→after pasting stress command→use top command

- In target group we can check whether the instances are healthy or not.



- After applying stress we have achieved max desired capacity.



- We can see that after achieving max desired capacity all instances generated are in healthy state .