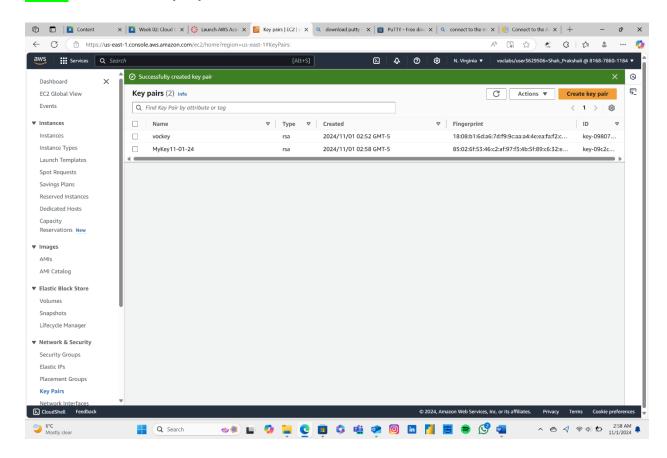# AWS EMR CLUSTER SETUP

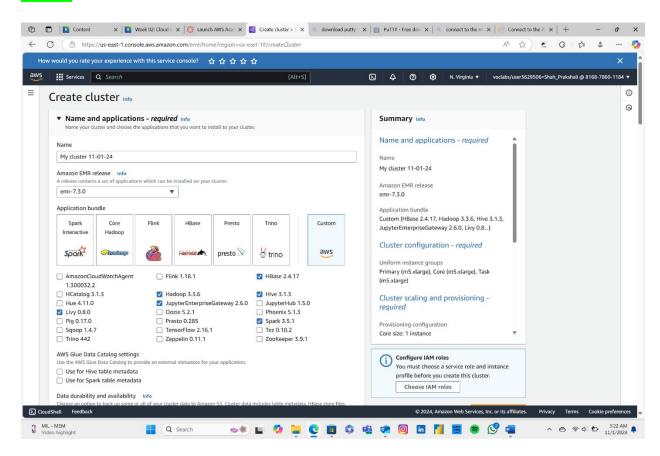Hnas-on Lab

- **Task 1:** Create a Security Key Pair

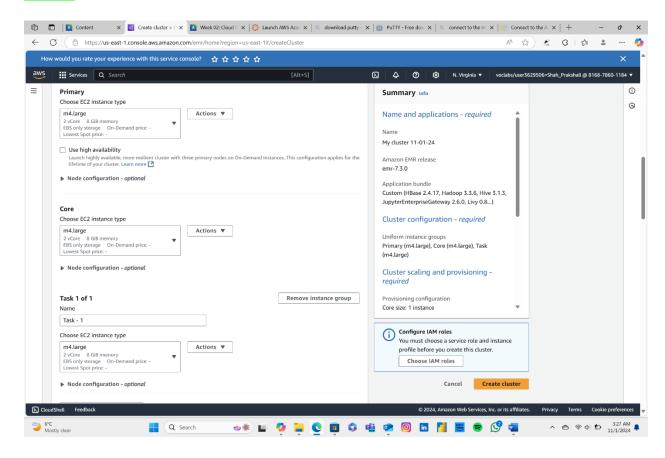==Task 2==

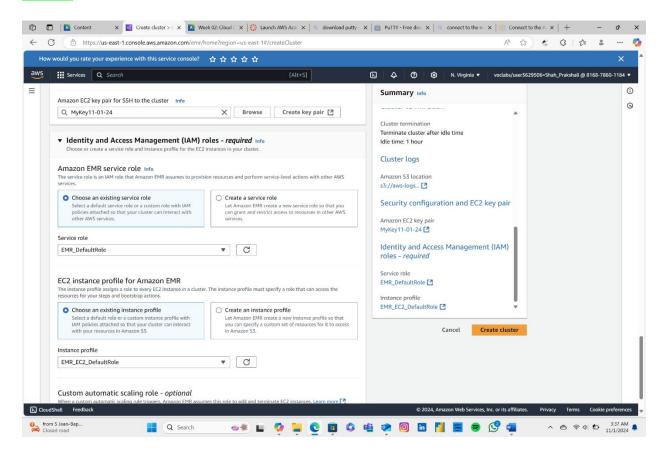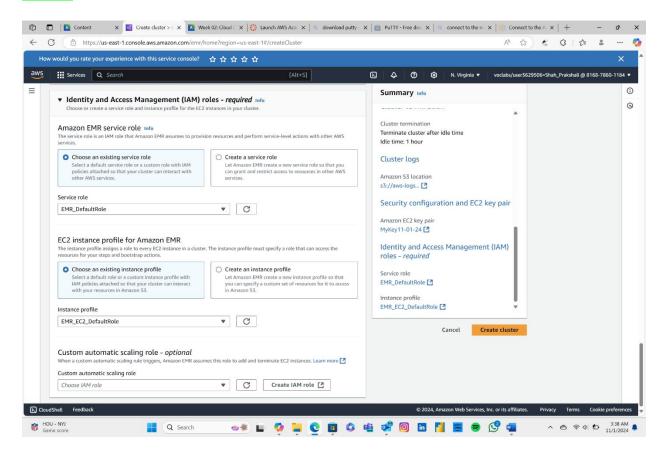- ==Task 2.1== Choose your software

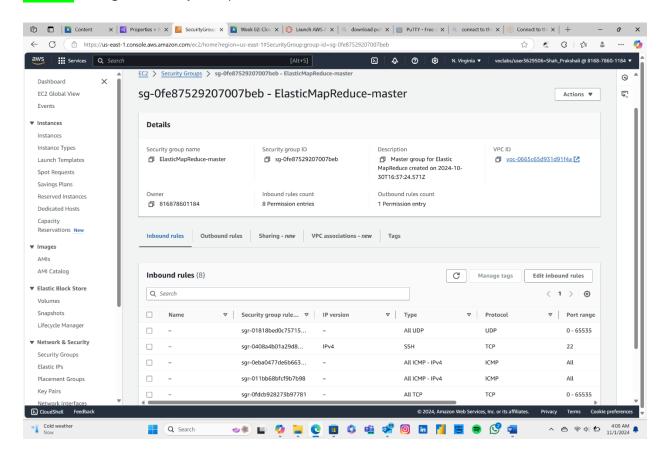- **Task 2.2** Configure Instance Group, instance type
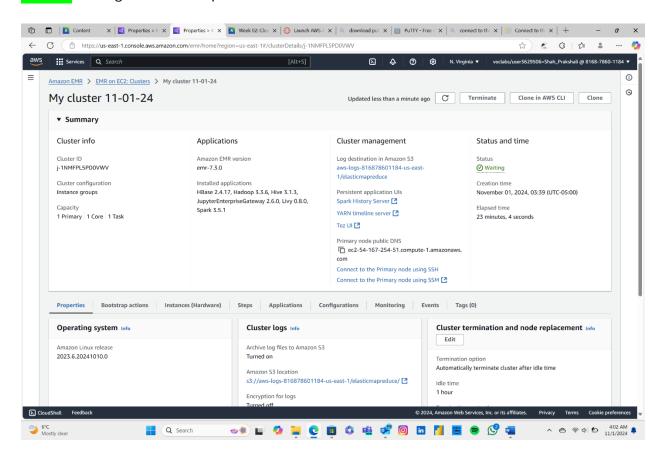
- **Task 2.3** Configure Security

- **Task 2.4** Configure IAM Role

- **Task 2.5** Configure Security Group Inbound Rules

- **Task 2.6** Waiting for EMR to be provisioned

- **Task 3** Connect to the EMR



# A. Reflection

Diving into the exercise of connecting to and logging into the Hadoop cluster was a comprehensive learning experience.

**What I Learned**
Configuring Security Key Pairs: I understood the significance of creating and managing security key pairs, which are essential for secure access to the EC2 instances.

Setting Up EMR Clusters: I learned how to navigate the AWS Management Console to set up an EMR cluster with specific software configurations like Hadoop, Spark, Hive, and HBase. Choosing the appropriate instance types (m4.large) was crucial for optimal performance.

Inbound Rules for SSH Access: Configuring the security group's inbound rules to allow SSH access from any IP address was a pivotal step, ensuring seamless connectivity.

**Challenges Faced**

Inbound Rules Configuration: One major challenge was ensuring the correct configuration of inbound rules to allow SSH access. Initially, I encountered issues with connecting to the EMR cluster due to misconfigured rules.

Instance Type Selection: Choosing the correct instance types for the primary, core, and task nodes also posed a challenge as selecting the wrong instance could impact performance.

Connectivity Issues: There were moments when connecting via Putty required rechecking the host name and ensuring the key was correctly saved and selected.

**Resolutions**

Referring to Instructions: Whenever I faced an issue, I revisited the detailed instructions and video guides provided in the task list. This ensured that I adhered to the steps meticulously.

Trial and Error: For configuring inbound rules, I used a trial-and-error approach, modifying settings until the SSH connection was successful.

Patience and Precision: Ensuring precision in following each step, from creating the security key pair to configuring the EMR cluster, played a vital role in overcoming challenges.

This exercise not only enhanced my technical skills in cloud-based infrastructure management but also underscored the importance of meticulous attention to detail and problem-solving in dynamic environments.