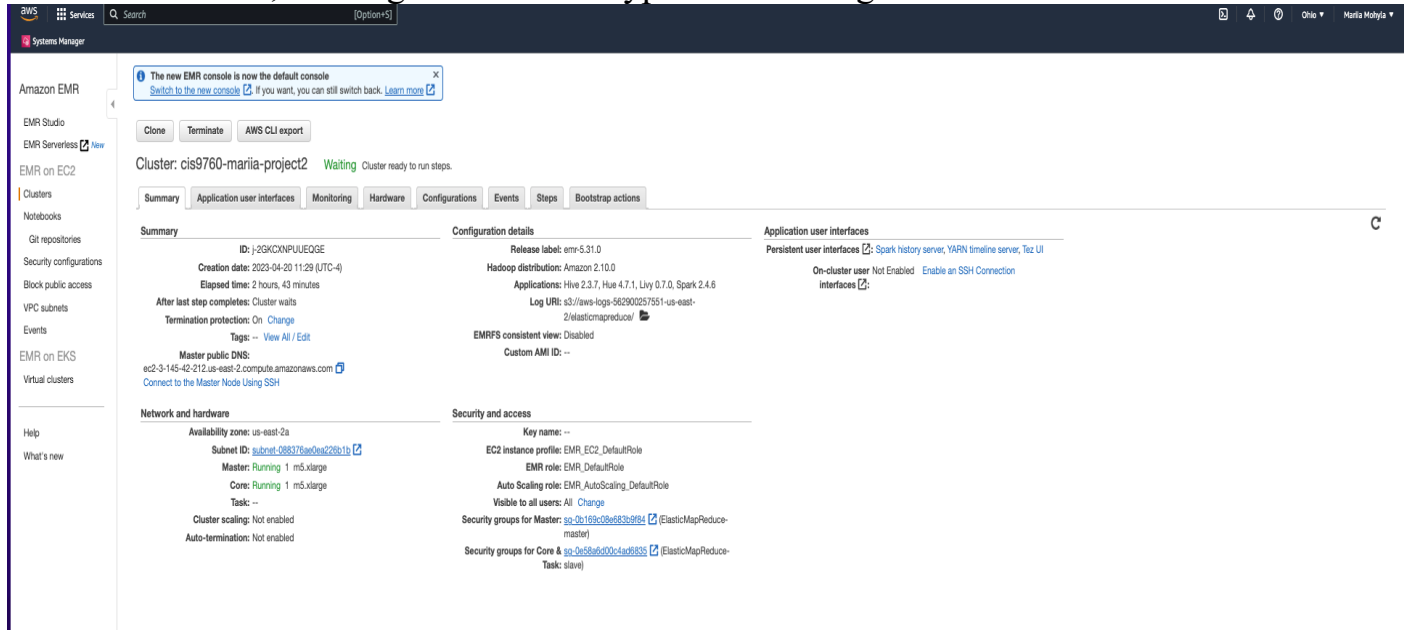The goal of this project is to provision a Spark infrastructure on AWS EMR ecosystem to load and run some exploratory data analysis using PySpark on IMDB's dataset from Kaggle located in publicly available S3 storage.
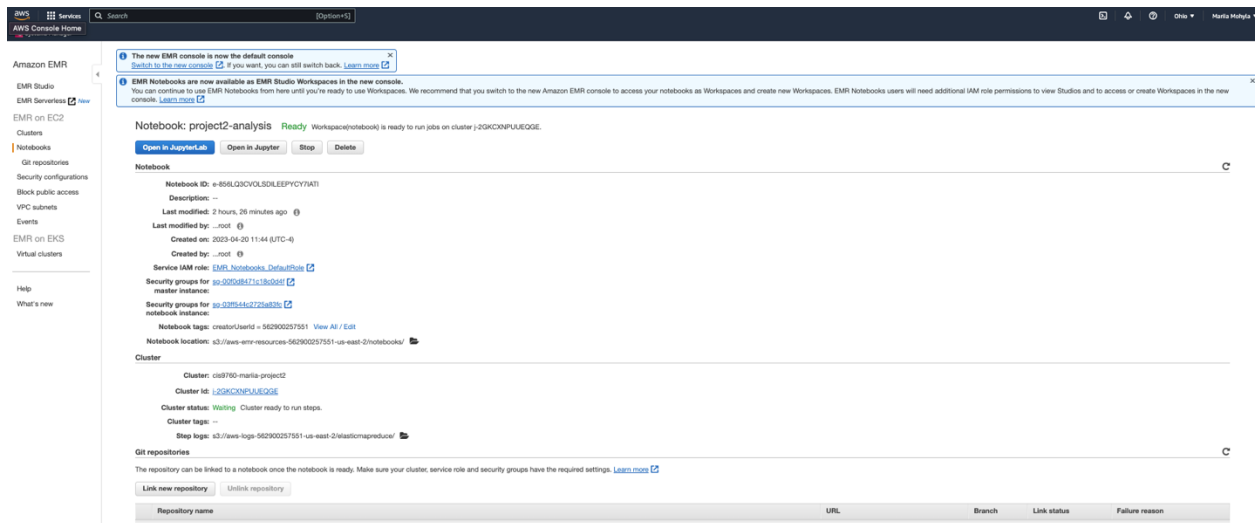
## STEP 1:
The project starts with setting up an EMR ecosystem to provision a Spark cluster of machines to run jobs of big data. For the project purposes, I chose 1 master node and 1 core node. Also, I configured instance types as m5.xlarge.



## STEP 2:
Then, I configured Jupyter Notebook to run on that cluster and enable the data analysis job. You may notice the configuration of the Jupyter Notebook below.

**STEP 3:**

I performed exploratory data analysis in my Jupyter Notebook by using PySpark library which allowed leveraging Dataframe API in Python and SQL queries to explore interesting facts about the IMDB dataset. The analysis consisted of four parts.

**Part I:** I started with installing some required libraries such as pandas and matplotlib and a general overview of the tables.

**Part II:** In the second part of the project, I discerned distinct movie genres and calculated the average rating per genre. Finally, I have plotted my findings. Short movies have the highest average rating while horror movies have the lowest.

**Part III:** Part 3 is related to the jobs in the movie industry. To be more specific, I built a bar chart with the top 5 job categories. Actor, actress, self, writer, and director are the top 5 jobs in the movie industry.

**Part IV:** Part 4 includes some additional questions that conclude my analysis.