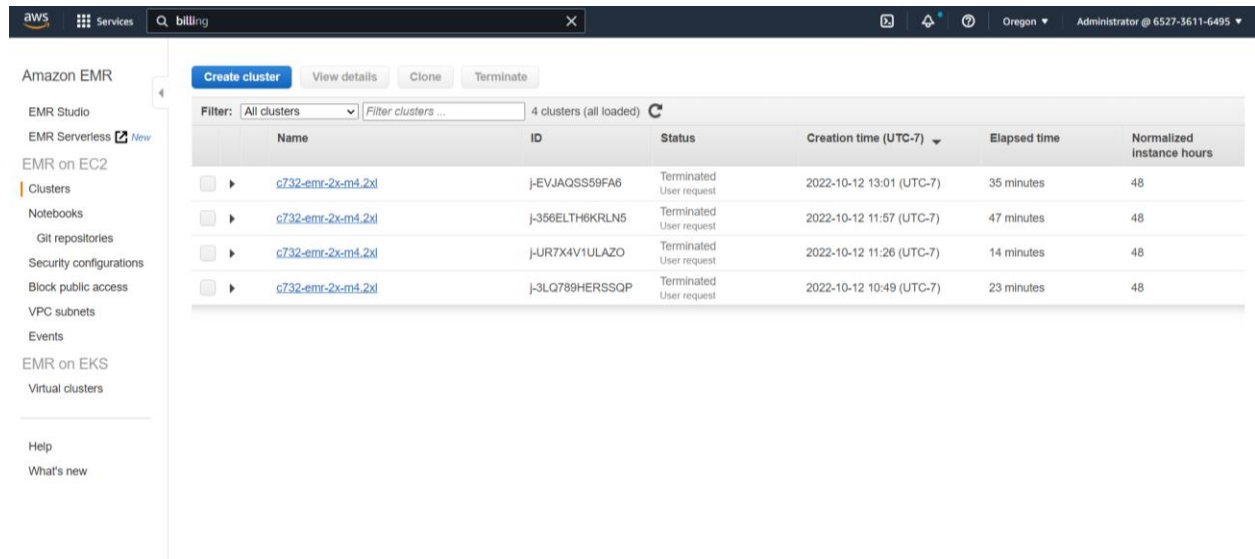


# ANSWERS

Q1.



The screenshot shows the Amazon EMR console interface. On the left is a navigation menu with options like EMR Studio, EMR Serverless, Clusters, Notebooks, etc. The main area displays a table of EMR clusters. The table has columns for Name, ID, Status, Creation time (UTC-7), Elapsed time, and Normalized instance hours. All four clusters listed are in a 'Terminated' status with a reason of 'User request'.

Name	ID	Status	Creation time (UTC-7)	Elapsed time	Normalized instance hours
<a href="#">c732-emr-2x-m4.2xl</a>	j-EVJAQSS59FA6	Terminated User request	2022-10-12 13:01 (UTC-7)	35 minutes	48
<a href="#">c732-emr-2x-m4.2xl</a>	j-356ELTH6KRLN5	Terminated User request	2022-10-12 11:57 (UTC-7)	47 minutes	48
<a href="#">c732-emr-2x-m4.2xl</a>	j-UR7X4V1ULAZO	Terminated User request	2022-10-12 11:26 (UTC-7)	14 minutes	48
<a href="#">c732-emr-2x-m4.2xl</a>	j-3LQ789HERSSQP	Terminated User request	2022-10-12 10:49 (UTC-7)	23 minutes	48

Figure1. Screenshot of all created cluster for assignment 5

Q2.

- Input file size without prefilter at S3 is 2.6 MiB and with filter it was 97.7 KiB. So, S3 filtered  $(2662 - 97.7)/2662$ , that is, 96.33%
- The filter, sorting (SQL where), file read options are done by both spark and S3.

Q3.

a)

The collect() job which included operations like reading text file and partitioning (Fig3) took the maximum time, making this application IO bound. The compute operations on the other hand took only 20-22s to complete. Even the last job, jobid5, took 1.3m which included writing the output to a file, again an IO task. (Fig4)

**m6gd.xlarge** costs \$0.1808 per hour, for 4 instances, (Memory = 16 GiB, Storage=1 x 237 NVMe SSD Network Storage=Up to 10 Gigabit). Therefore, the total costs would be  $0.1808 * 4 * \text{hours the cluster was running}$ . If dataset reddit-5 was 10 times larger the costs would be 10 times more. This can be avoided if we can split the large reddit-5 data into smaller more manageable files to process on 16 instances.