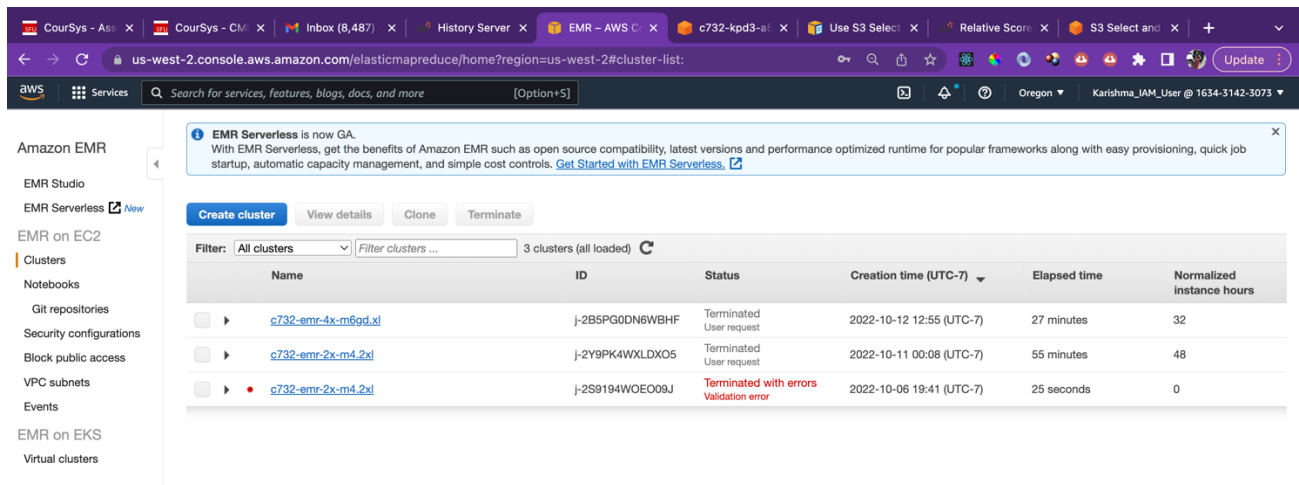


Q1) The terminated clusters are as shown below.



Amazon EMR console screenshot showing a list of clusters. The table displays the following information:

Name	ID	Status	Creation time (UTC-7)	Elapsed time	Normalized instance hours
<a href="#">c732-emr-4x-m6gd.xl</a>	j-2B5PG0DN6WBHF	Terminated User request	2022-10-12 12:55 (UTC-7)	27 minutes	32
<a href="#">c732-emr-2x-m4.2xl</a>	j-2Y9PK4WXLDO5	Terminated User request	2022-10-11 00:08 (UTC-7)	55 minutes	48
<a href="#">c732-emr-2x-m4.2xl</a>	j-2S9194W0EO09J	Terminated with errors Validation error	2022-10-06 19:41 (UTC-7)	25 seconds	0

Q2) a) The input size without filtering = 2.6 MiB

The input size with S3 filtering = 97.7 KiB

The fraction of the input file was filtered and sent to Spark =  $97.7/2662.4 = 0.0366$

Thus the amount that was filtered out by S3Select = 96.34% and the amount that was sent to spark after filtering = 3.66%

### Details for Stage 0 (Attempt 0)

Resource Profile Id: 0  
 Total Time Across All Tasks: 14 s  
 Locality Level Summary: Rack local: 4  
 Input Size / Records: 2.6 MiB / 3245  
 Output Size / Records: 27.2 KiB / 3245  
 Associated Job Ids: 0

Without S3 Select

### Details for Stage 0 (Attempt 0)

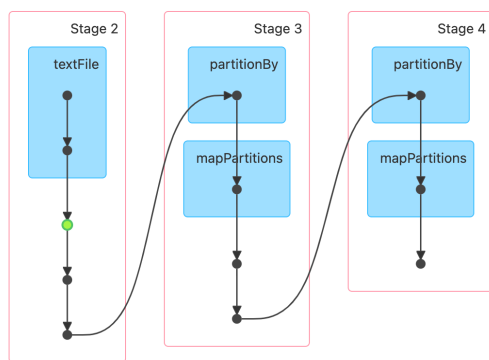
Resource Profile Id: 0  
 Total Time Across All Tasks: 10 s  
 Locality Level Summary: Rack local: 4  
 Input Size / Records: 97.7 KiB / 3245  
 Output Size / Records: 27.2 KiB / 3245  
 Associated Job Ids: 0

With S3 Select

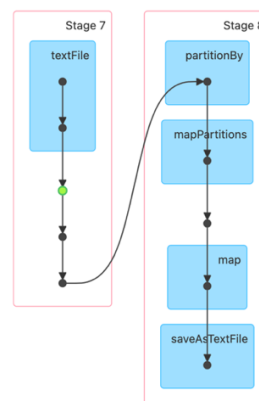
b) The SQL operations such as filtering, where and select, etc are being performed through S3Select whereas the other operations are performed by Spark.

Q3) From the Jobs and DAGs below we can see that the maximal amount of time is taken through the collect job and the runJob, both involve reading from the input and writing into the output, respectively. Since most time is spent in IO operations compared to the rest of the computations, we can say that the application is IO bound.

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
5	runJob at SparkHadoopWriter.scala:83 runJob at SparkHadoopWriter.scala:83	2022/10/12 20:14:32	1.9 min	2/2	32/32
4	sortBy at /mnt/tmp/spark-5b07d63a-4184-433c-a192-45bc1483b0a8/relative_score_bcast.py:59 sortBy at /mnt/tmp/spark-5b07d63a-4184-433c-a192-45bc1483b0a8/relative_score_bcast.py:59	2022/10/12 20:13:50	41 s	1/1	16/16
3	sortBy at /mnt/tmp/spark-5b07d63a-4184-433c-a192-45bc1483b0a8/relative_score_bcast.py:59 sortBy at /mnt/tmp/spark-5b07d63a-4184-433c-a192-45bc1483b0a8/relative_score_bcast.py:59	2022/10/12 20:13:11	39 s	1/1	16/16
2	collect at /mnt/tmp/spark-5b07d63a-4184-433c-a192-45bc1483b0a8/relative_score_bcast.py:52 collect at /mnt/tmp/spark-5b07d63a-4184-433c-a192-45bc1483b0a8/relative_score_bcast.py:52	2022/10/12 20:10:31	2.7 min	3/3	48/48
1	sortBy at /mnt/tmp/spark-5b07d63a-4184-433c-a192-45bc1483b0a8/relative_score_bcast.py:48 sortBy at /mnt/tmp/spark-5b07d63a-4184-433c-a192-45bc1483b0a8/relative_score_bcast.py:48	2022/10/12 20:09:56	35 s	1/1	16/16
0	sortBy at /mnt/tmp/spark-5b07d63a-4184-433c-a192-45bc1483b0a8/relative_score_bcast.py:48 sortBy at /mnt/tmp/spark-5b07d63a-4184-433c-a192-45bc1483b0a8/relative_score_bcast.py:48	2022/10/12 20:07:13	2.7 min	1/1	16/16



Collect DAG



runJob DAG

b)

cost of m6gd.xlarge = \$.01808/hour of usage while using 4vCPUs

Since the cost is calculated depending on the time taken, a dataset 10 times the size of reddit-5 would take much more time to execute with the same 4 instances and hence the price would increase accordingly.

In order to process a very large dataset while making use of 16 instances, we could partition the dataset so that the instances can work parallelly with these chunks thus making better use of the available resources.