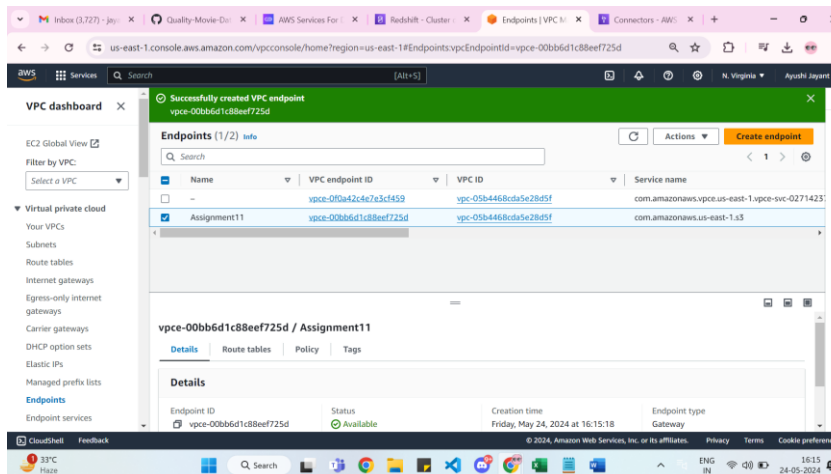
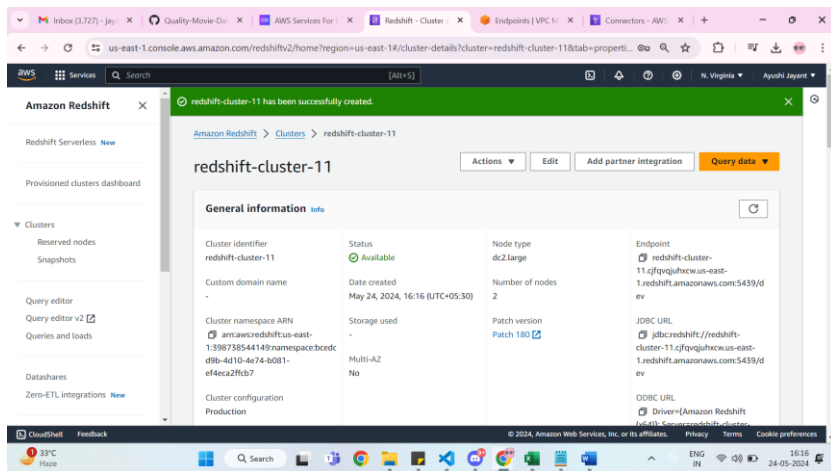
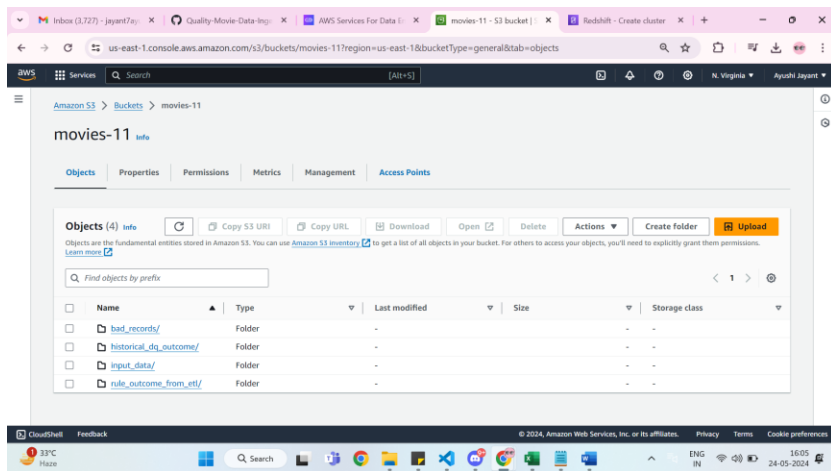
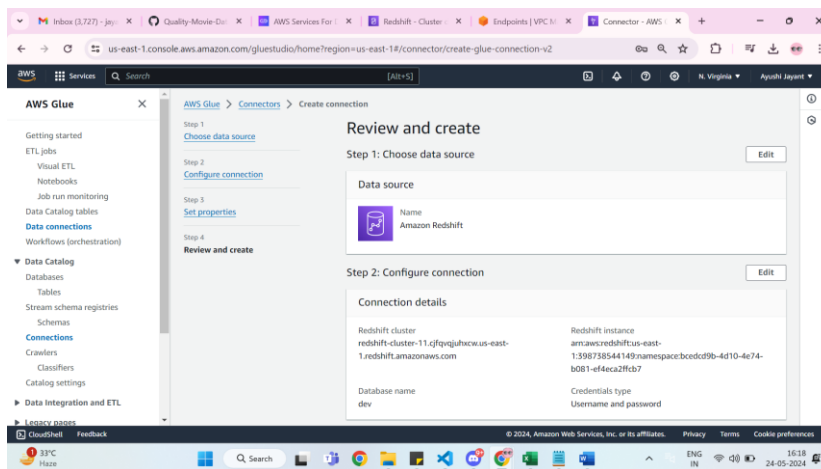
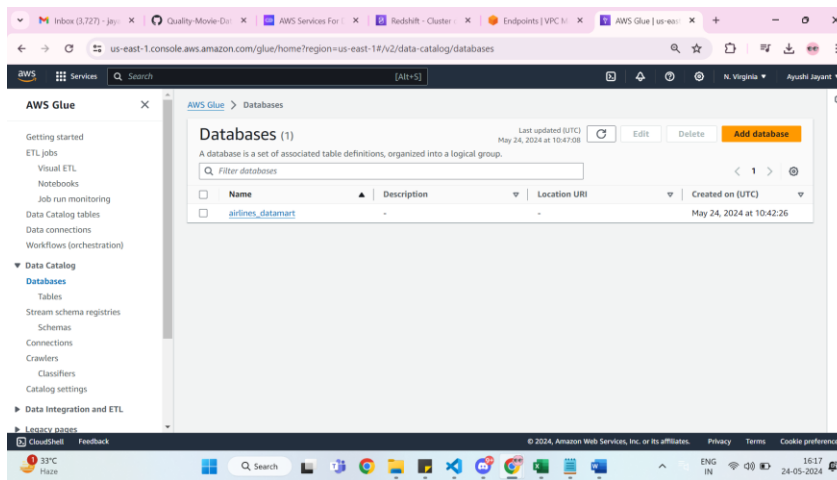


Assignment11:Movie Quality Data Ingestion

- 1)We have created relevant folders in s3 bucket with movie data inside input_data folder in s3 bucket
- 2)created redshift cluster
- 3)created vpc security groups inbound and outbound rules along with s3 bucket as vpc endpoint
- 4)created glue database as named airlines_datamart
- 5)created redshift schema and table for movie data in query editor v2
- 6)we have made a data quality rule set
- 7)created visual etl glue job and set the eventbridge rule according to the data quality rule
- 8)After successfully running glue job we check the table in redshift with the relevant columns
- 9)We have also received bad records in our folder located in s3 bucket
- 10)created a materialised view from our redshift table and query on basis of that
- 11)created step function to orchestrate our pipeline and attached sns topic with the same for success/failure notifications
- 12)We have created an eventbridge rule to invoke the step function as soon as an object is created in the movie-data bucket
- 13)step functions ran successfully and we got the notification on mail
- 14)After querying in the redshift, we get the total count of records as 99 from 33 because we have the refreshed data in the materialised view





Create schema movies;

CREATE TABLE movies.imdb_movies_rating (

Poster_Link VARCHAR(MAX),

Series_Title VARCHAR(MAX),

Released_Year VARCHAR(10),

Certificate VARCHAR(50),

Runtime VARCHAR(50),

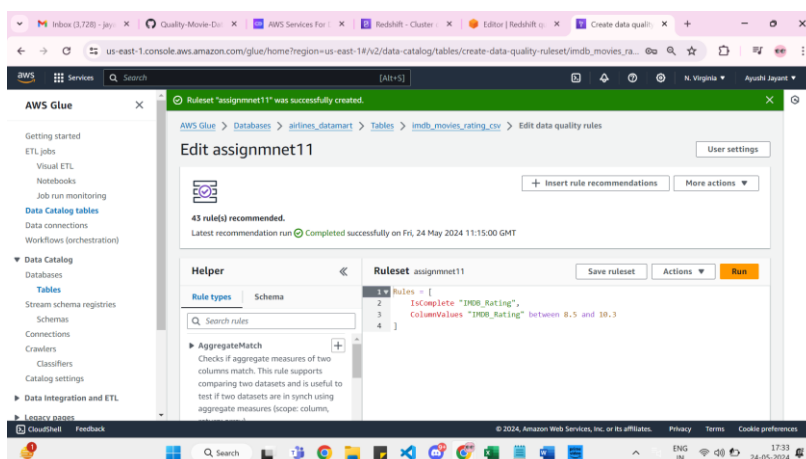
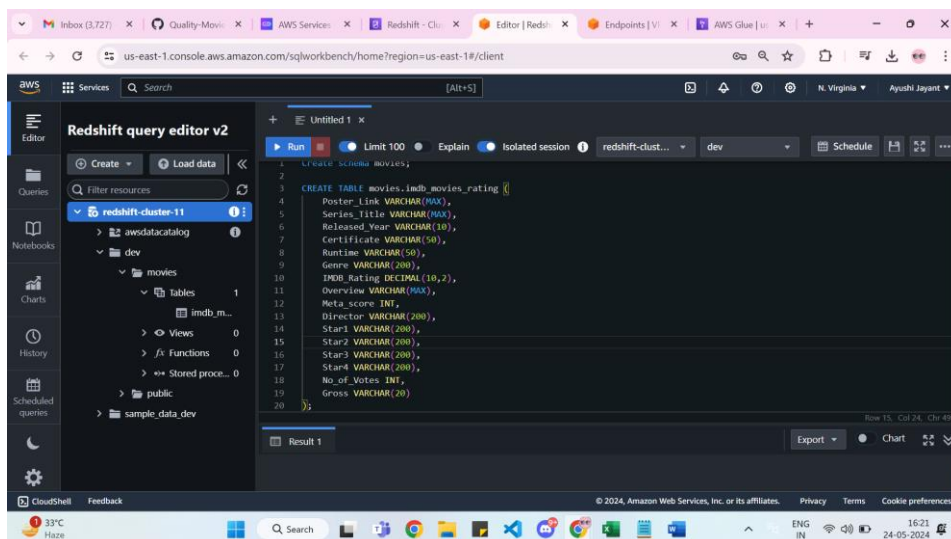
Genre VARCHAR(200),

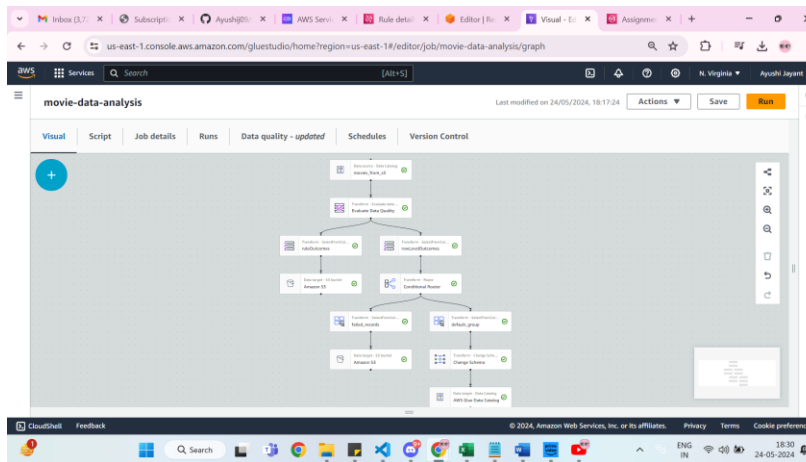
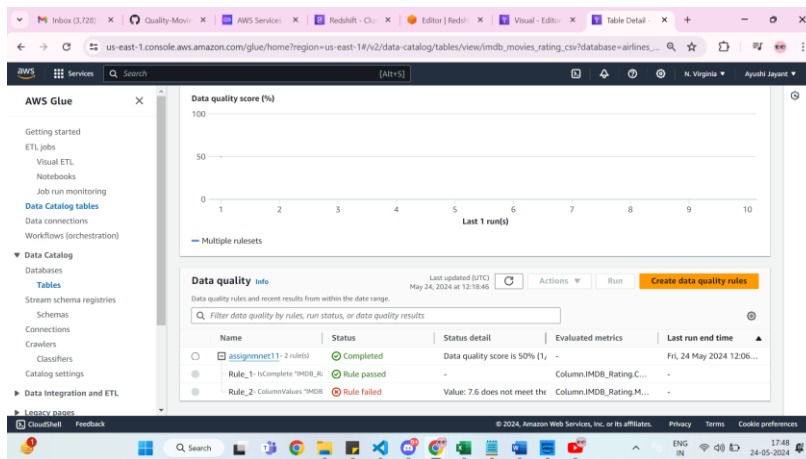
IMDB_Rating DECIMAL(10,2),

Overview VARCHAR(MAX),

Meta_score INT,

Director VARCHAR(200),
Star1 VARCHAR(200),
Star2 VARCHAR(200),
Star3 VARCHAR(200),
Star4 VARCHAR(200),
No_of_Votes INT,
Gross VARCHAR(20)
);





us-east-1.console.aws.amazon.com/events/home?region=us-east-1#/eventbus/default/rules/movie-data-quality11

movie-data-quality11

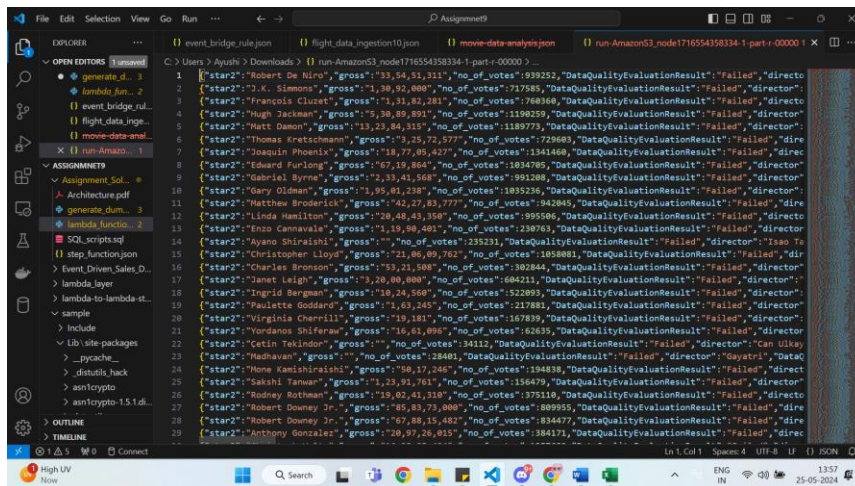
Edit Disable Delete CloudFormation Template

Rule details

Rule name	Status	Event bus name	Type
movie-data-quality11	Enabled	default	Standard

Event pattern

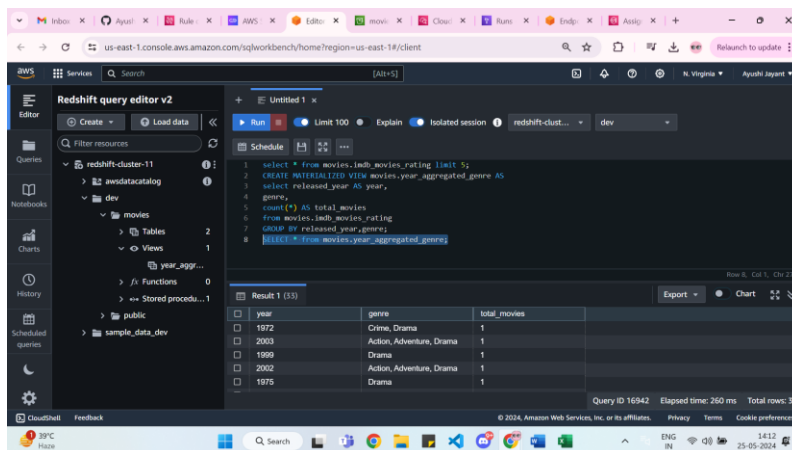
```
1 {
2   "source": ["aws-glue-dataquality"],
3   "detail-type": ["Data Quality Evaluation Results Available"],
4   "detail": {
5     "status": ["SUCCESS", "FAILURE"]
6   }
7 }
```

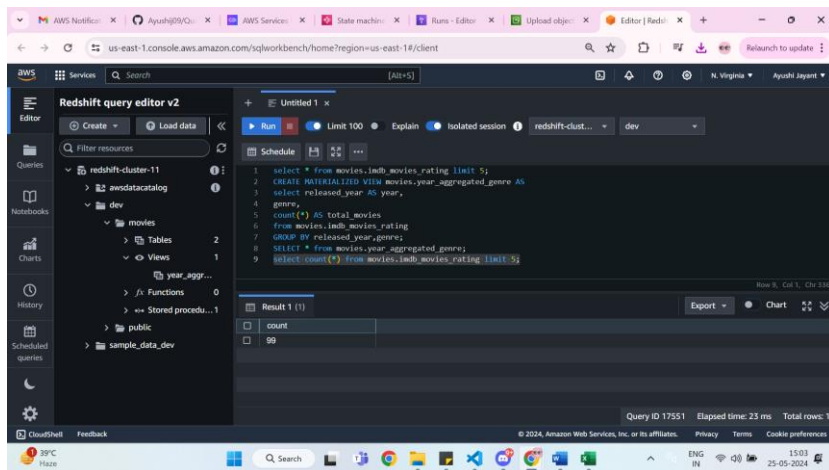
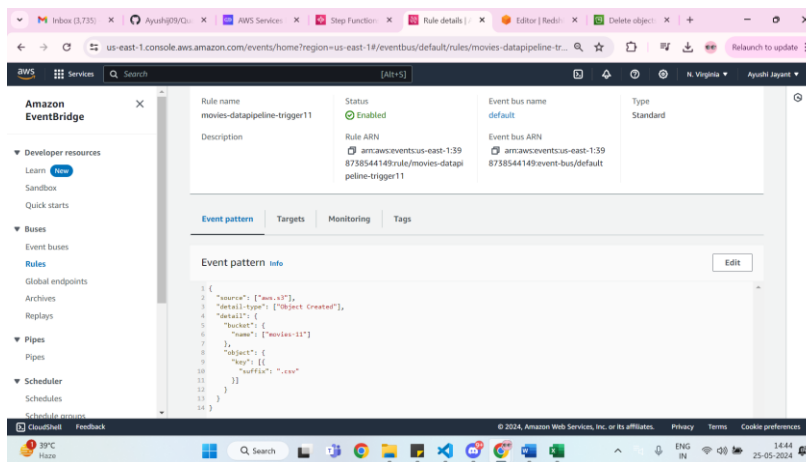



```

CREATE MATERIALIZED VIEW movies.year_aggregated_genre AS
select released_year AS year,
genre,
count(*) AS total_movies
from movies.imdb_movies_rating
GROUP BY released_year,genre;

```





Now, if we check after refreshing the materialised view, the count has been tripled i.e. from 33→99

us-east-1.console.aws.amazon.com/iglworkebench/home?region=us-east-1#/client

Redshift query editor v2

Filter resources

- redshift-cluster-11
 - awsdatacatalog
 - dev
 - movies
 - tbl_movies_rating 2
 - tbl_year_aggreg... 1
 - Views
 - year_aggregated_gene 1
 - Functions 0
 - Stored procedures 1
 - public
 - sample_data_dev

```
1 select * from movies.tbl_movies_rating limit 5;
2 CREATE MATERIALIZED VIEW movies.year_aggregated_genres AS
3 select released_year, AS year,
4 genres,
5 count(*) AS total_movies
6 from movies.tbl_movies_rating
7 GROUP BY released_year, genres;
8 SELECT * from movies.year_aggregated_genres;
9 select count(*) from movies.tbl_movies_rating limit 5;
10 refresh MATERIALIZED VIEW movies.year_aggregated_genres;
```

Result 1 (53)

year	genre	total_movies
1980	Action, Adventure, Fantasy	3
1983	Biography, Drama, History	3
1987	Crime, Drama	3
1994	Drama, Romance	3

Query ID: 17648 Elapsed time: 8971 ms Total rows: 35