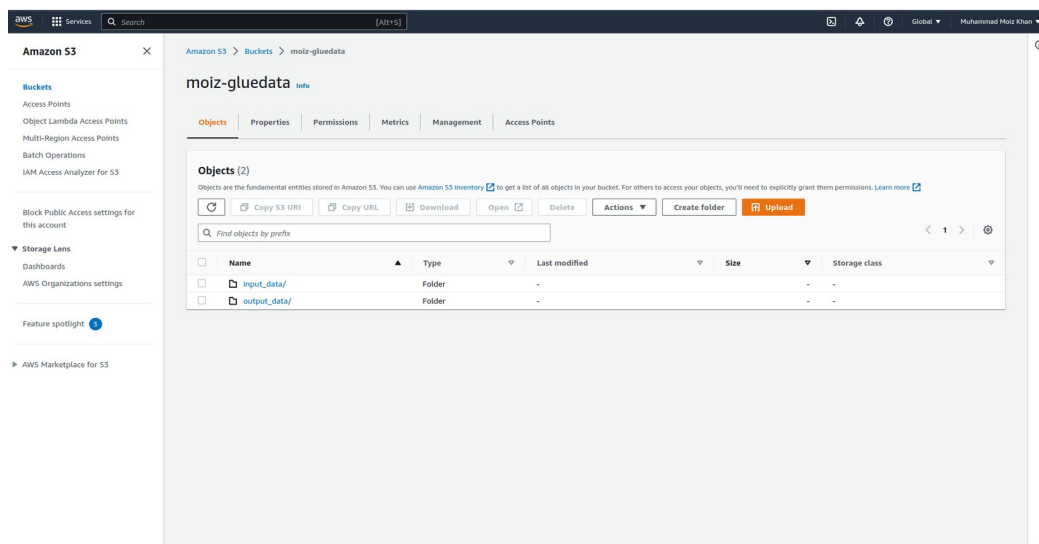


Talha Khan (2303.009.KHI.DEG)
Muhammad Moiz Khan (2303.022.KHI.DEG)

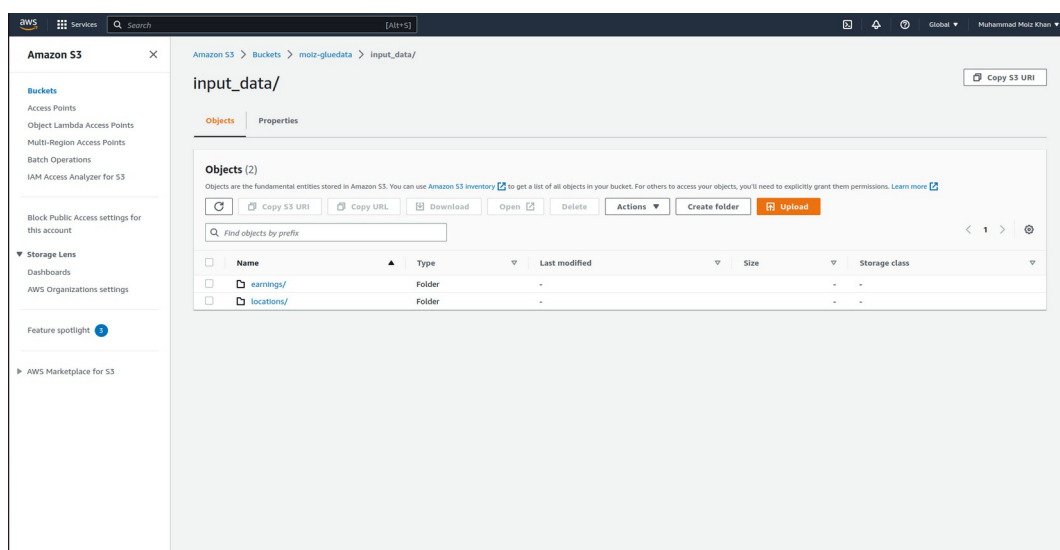
ASSIGNMENT 5.2

Using the salary CSV as a base, prepare a new data file with employees' office locations. Make sure there are 5-6 distinct locations that are shared between employees. Create a Glue job that aggregates the data based on the office location to calculate average salaries and raise percentages for these locations.

1) First create a folder in s3 bucket name (input and output data):



2) Next, we store a two dataset earning.csv and location.csv in input folder:



3) Now we create two crawlers named earnings and locations respectively, and extract the meta data :

The screenshot displays the AWS Glue console interface. A green banner at the top indicates "Crawler successfully starting" for "moiz_s3_earnings_crawler". The left sidebar shows the navigation menu with "Crawlers" selected under the "Data Catalog" section. The main content area shows the configuration for the "moiz_s3_locations_crawler".

Crawler properties

Property	Value
Name	moiz_s3_locations_crawler
IAM role	talhakhah-glue-role
Database	talhakhah_glue_database
State	READY
Description	-
Security configuration	-
Lake Formation configuration	-
Table prefix	-
Maximum table threshold	-

Crawler runs (1)

Start time (UTC)	End time (UTC)	Current/last duration	Status	DPU hours	Table changes
May 22, 2023 at 05:05:44	May 22, 2023 at 05:06:28	44 s	Completed	0.072	1 table change, 0 partition changes

The screenshot displays the AWS Glue console interface for the "moiz_s3_earnings_crawler". A green banner at the top indicates "Crawler successfully starting" for "moiz_s3_earnings_crawler". The left sidebar shows the navigation menu with "Crawlers" selected under the "Data Catalog" section. The main content area shows the configuration for the "moiz_s3_earnings_crawler".

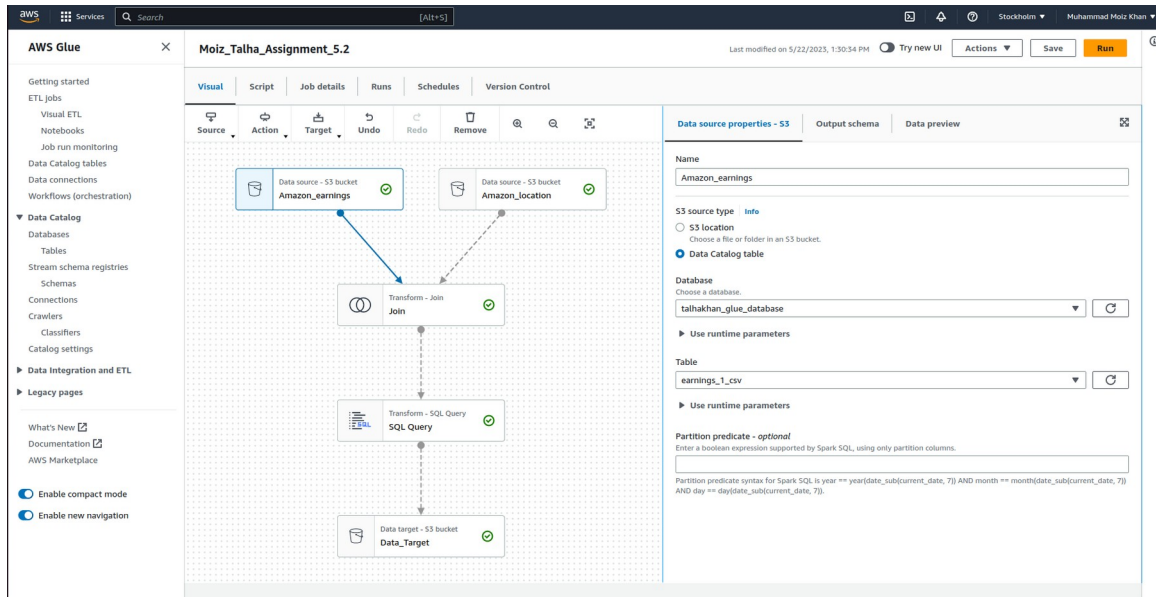
Crawler properties

Property	Value
Name	moiz_s3_earnings_crawler
IAM role	talhakhah-glue-role
Database	talhakhah_glue_database
State	READY
Description	-
Security configuration	-
Lake Formation configuration	-
Table prefix	moiz_
Maximum table threshold	-

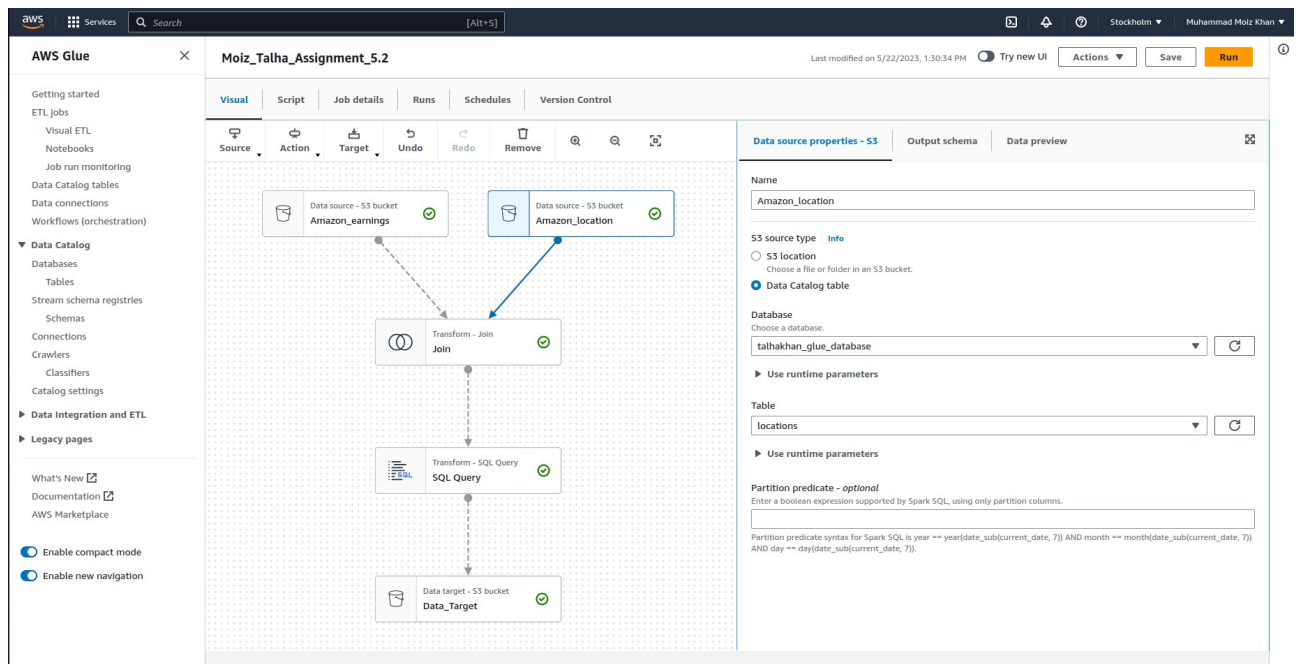
Crawler runs (1)

Start time (UTC)	End time (UTC)	Current/last duration	Status	DPU hours	Table changes
May 22, 2023 at 07:35:32	May 22, 2023 at 07:36:20	47 s	Completed	0.066	1 table change, 0 partition changes

4) After creating the glue crawler we start creating the job:



5) Now we create two s3 source one is employee earning data and location, we perform inner join on both data on emp_id and after that prepare for querying :



6) Now finally we load the data and show in the output folder :

The screenshot shows a web-based data transformation tool interface. At the top, there's a header with a user profile 'Muhammad Moiz Khan' and location 'Stockholm'. Below the header, a status bar indicates 'Last modified on 5/22/2023, 1:42:44 PM' and includes buttons for 'End session', 'Actions', 'Save', and 'Run'. A green banner with a close button is visible below the status bar. The main workspace has three tabs: 'Transform', 'Output schema', and 'Data preview'. The 'Data preview' tab is active, showing a table with 3 columns: 'location', 'average_earnings', and 'raise_percentage'. The table contains 3 rows of data. A search bar with the placeholder 'Filter sample dataset' is located above the table. The table data is as follows:

location	average_earnings	raise_percentage
B	6286.75	155.14407467532467
A	5926.05	191.49286768322676
C	5257	39.29517753047165

The screenshot shows the AWS S3 console interface. The left sidebar contains navigation options like 'Buckets', 'Access Points', 'Object Lambda Access Points', 'Multi-Region Access Points', 'Batch Operations', 'IAM Access Analyzer for S3', 'Storage Lens', 'Dashboards', 'AWS Organizations settings', and 'Feature spotlight'. The main content area shows the path 'Amazon S3 > Buckets > moiz-gluedata > output_data/ > earnings_with_folder/'. The 'Objects' tab is selected, displaying a list of 5 objects. The objects are all parquet files, each 599.0 B in size, and are stored in the 'Standard' storage class. The objects are named as follows:

- run-1684744292621-part-block-0-r-00002-snappy.parquet
- run-1684744292621-part-block-0-r-00014-snappy.parquet
- run-1684744292621-part-block-0-r-00021-snappy.parquet
- run-1684744292621-part-block-0-r-00025-snappy.parquet
- run-1684744292621-part-block-0-r-00031-snappy.parquet

