# Unit 5.2 Graded Assignment:
## Syed Muhammad Raqim Ali Shah (2303.KHI.DEG.008)
## Maaz Javaid Siddique (2303.KHI.DEG.004)
## Qadeer Hussain (2303.KHI.DEG.006)
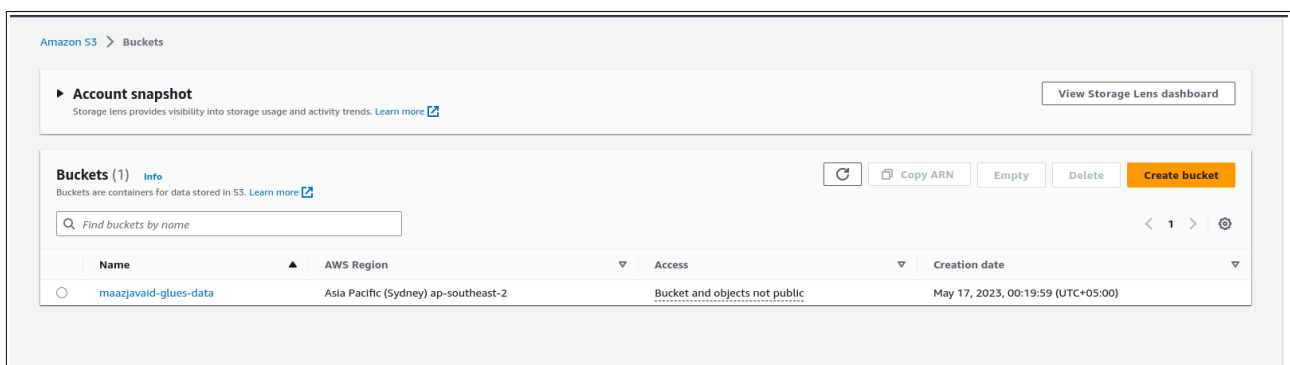
## Daily Assignment :

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.
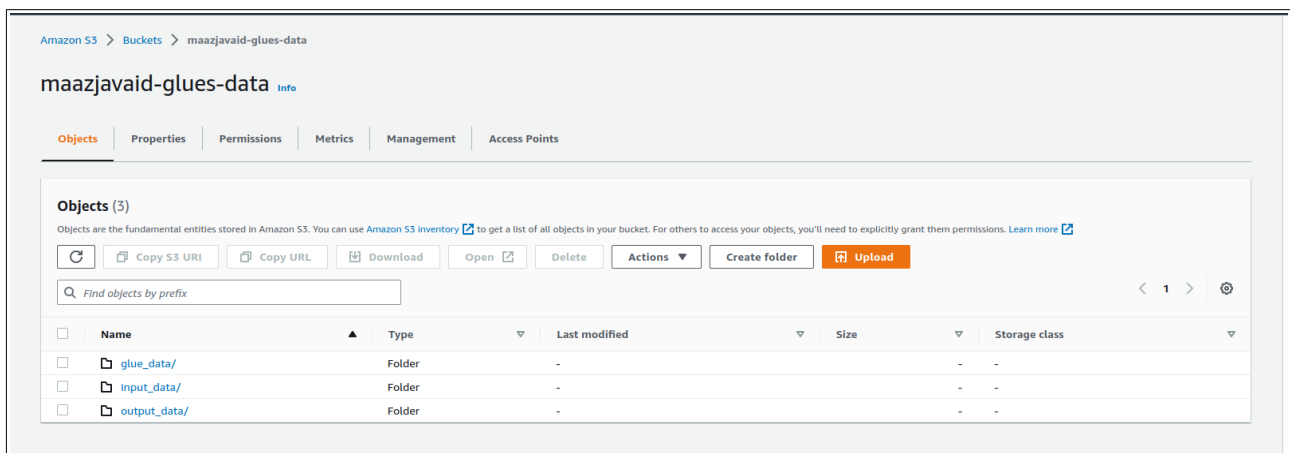
## Answer:

### 1.Set up an AWS Glue job in the AWS Glue console

First of all we create the S3 bucket:



Then we create the directories on S3 bucket:

# input_data/

Copy S3 URI

**Objects** | Properties

## Objects (2)

Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ↗ to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more ↗

Copy S3 URI | Copy URL | Download | Open ↗ | Delete | Actions ▼ | Create folder | **Upload**

Find objects by prefix

| | Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|---|
| ☐ | 📄 earnings/ | Folder | - | - | - |
| ☐ | 📄 location/ | Folder | - | - | - |

---

# output_data/

Copy S3 URI

**Objects** | Properties

## Objects (2)

Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ↗ to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more ↗

Copy S3 URI | Copy URL | Download | Open ↗ | Delete | Actions ▼ | Create folder | **Upload**

Find objects by prefix

| | Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|---|
| ☐ | 📄 employees_earnings/ | Folder | - | - | - |
| ☐ | 📄 location/ | Folder | - | - | - |

Create the RDS database:

RDS > Databases

ⓘ **Consider creating a Blue/Green Deployment to minimize downtime during upgrades**
You may want to consider using Amazon RDS Blue/Green Deployments and minimize your downtime during upgrades. A Blue/Green Deployment provides a staging environment for changes to production databases. RDS User Guide ↗ Aurora User Guide ↗

## Databases (1)

⬤ Group resources | Modify | Actions ▼ | Restore from S3 | **Create database**

Filter by databases

| | DB identifier ▲ | Role ▽ | Engine ▽ | Status ▽ | Region & AZ ▽ | Size ▽ | Actions ▽ | CPU ▽ | Current activity ▽ | Maintenance ▽ | VPC |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ○ | maazjavaid-employees-db | Instance | PostgreSQL | ✓ Available | ap-southeast-2b | db.t3.micro | 5 Actions | 3.83% | 0 Connections | none | vpc-091369db5 |

Then we go to IAM service and create the IAM Role:

## Create the VPC Endpoint:



## Security group rules:

# 2.ETL in Glue

## Create database:



## Glue Crawlers:



## Create job:

Visual | Script | Job details | Runs | Schedules | Version Control

Source | Action | Target | Undo | Redo | Remove

**Data source properties - S3** | Output schema | Data preview

Data source – S3 bucket
Amazon S3

Data source – S3 bucket
Amazon S3

Transform - Join
Join

Transform - SQL Query
SQL Query

Data target – S3 bucket
Amazon S3

**Name**
Amazon S3

**S3 source type**   Info
○ S3 location
    Choose a file or folder in an S3 bucket.
● Data Catalog table

**Database**
Choose a database.
maazjavaid_glue_database

▶ Use runtime parameters

**Table**
maazjavaid_earnings

▶ Use runtime parameters

**Partition predicate – optional**
Enter a boolean expression supported by Spark SQL, using only partition columns.

Partition predicate syntax for Spark SQL is year == year(date_sub(current_date, 7)) AND month == month(date_sub(current_date, 7)) AND day == day(date_sub(current_date, 7)).

---

**Data source properties - S3** | Output schema | Data preview

Data source – S3 bucket
Amazon S3

Data source – S3 bucket
Amazon S3

Transform - Join
Join

Transform - SQL Query
SQL Query

Data target – S3 bucket
Amazon S3

**Name**
Amazon S3

**S3 source type**   Info
○ S3 location
    Choose a file or folder in an S3 bucket.
● Data Catalog table

**Database**
Choose a database.
maazjavaid_glue_database

▶ Use runtime parameters

**Table**
maazjavaid_location

▶ Use runtime parameters

**Partition predicate – optional**
Enter a boolean expression supported by Spark SQL, using only partition columns.

Partition predicate syntax for Spark SQL is year == year(date_sub(current_date, 7)) AND month == month(date_sub(current_date, 7)) AND day == day(date_sub(current_date, 7)).

---

**Transform** | Output schema | Data preview

Data source – S3 bucket
Amazon S3

Data source – S3 bucket
Amazon S3

Transform - Join
Join

Transform - SQL Query
SQL Query

Data target – S3 bucket
Amazon S3

**Name**
Join

**Node parents**
Choose which nodes will provide inputs for this one.
Choose one or more parent nodes

Amazon S3 ✕   Amazon S3 ✕
S3 – DataSource   S3 – DataSource

ⓘ The parents of this node have overlapping field names. AWS Glue Studio can add an Apply Mapping node to rename them and avoid downstream issues.

**Custom prefix**
Add a prefix to the field names of the parent node on the right
right   Resolve it

**Join type**
Select the type of join to perform.
Inner join
Select all rows from both datasets that meet the join condition.

**Join conditions**
Select a field from each parent node for the join condition.

Amazon S3      Amazon S3
emp_id   =   emp_id

Add condition

**maazjavaid_employee_earnings_job**

Last modified on 5/17/2023, 12:12:04 PM · Try new UI · Actions ▼ · Save · Run

Visual | Script | Job details | Runs | Schedules | Version Control

Source | Action | Target | Undo | Redo | Remove

Data source - S3 bucket — Amazon S3 ✓
Data source - S3 bucket — Amazon S3 ✓
Transform - Join — Join
Transform - SQL Query — SQL Query
Data target - S3 bucket — Amazon S3

Transform | Output schema | Data preview

**Name**
SQL Query

**Node parents**
Choose which nodes will provide inputs for this one.
Choose one or more parent nodes ▼
Join — Join - Transform ✕

**Associate an alias with each input source**  Info
Edit the aliases used for the inputs to this node.

| Input sources | SQL aliases |
|---|---|
| Join | myDataSource |

**SQL query**
Enter a SQL statement to add to your job.

```
1  SELECT
2      location,
3      AVG(earnings) AS average_earning,
4      (AVG(earnings)-MIN(earnings))/MIN(earnings)*100 AS raise_percentage
5  FROM
6      myDataSource
7  GROUP BY
8      location;
9  |
```



**maazjavaid_employee_earnings_job**

Last modified on 5/17/2023, 12:12:04 PM · Try new UI · Actions ▼ · Save · Run

Visual | Script | Job details | Runs | Schedules | Version Control

Source | Action | Target | Undo | Redo | Remove

Data source - S3 bucket — Amazon S3 ✓
Data source - S3 bucket — Amazon S3 ✓
Transform - Join — Join
Transform - SQL Query — SQL Query
Data target - S3 bucket — Amazon S3

Data target properties - S3 | Output schema | Data preview

**Name**
Amazon S3

**Node parents**
Choose which nodes will provide inputs for this one.
Choose one or more parent nodes ▼
SQL Query — SqlCode - Transform ✕

**Format**
Parquet ▼

**Compression Type**
Snappy ▼

**S3 Target Location**
Choose an S3 location in the format s3://bucket/prefix/object/ with a trailing slash (/).
s3://maazjavaid-glues-data/output_data/location/  ✕  View ⧉  Browse S3

**Data Catalog update options**  Info
Choose how you want to update the Data Catalog table's schema and partitions. These options will only apply if the Data Catalog table is an S3 backed source.
○ Do not update the Data Catalog
● Create a table in the Data Catalog and on subsequent runs, update the schema and add new partitions
○ Create a table in the Data Catalog and on subsequent runs, keep existing schema and add new partitions

**Database**
Choose the database from the AWS Glue Data Catalog.
maazjavaid_glue_database ▼  ↻

▶ Use runtime parameters



**maazjavaid_employee_earnings_job**

Last modified on 5/17/2023, 12:49:36 PM · Try new UI · End session · Actions ▼ · Save · Run

Visual | Script | Job details | Runs | Schedules | Version Control

Source | Action | Target | Undo | Redo | Remove

Data source - S3 bucket — Amazon S3 ✓
Data source - S3 bucket — Amazon S3 ✓
Transform - Join — Join ✓
Transform - SQL Query — SQL Query
Data target - S3 bucket — Amazon S3 ✓

Transform | Output schema | Data preview

**Data preview** (5)  Info                    Previewing 3 of 3 fields

Filter sample dataset

| location | average_earning | raise_percentage |
|---|---|---|
| B | 6086.875 | 184.30056048575432 |
| C | 5695.3 | 158.9949977262392 |
| A | 6217.975 | 205.85218888342354 |
| D | 5635.075 | 180.91101694915253 |
| E | 5503.4 | 154.31608133086874 |

Successfully updated job: