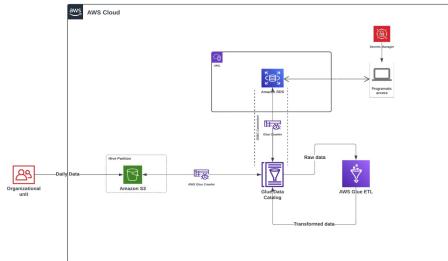


PS/Objective:

Implement an end-to-end data pipeline on AWS, filtering prepaid telecom data from daily CSV uploads to S3 using Glue ETL with distributed computation and securely ingesting them into RDS. Covering RDS setup, Glue catalog creation, JDBC connection establishment, visual ETL job creation with job bookmarking, and distributed computation using Glue ETL.



Step 1: Create RDS with the Following Parameters

1. Select RDS
 2. Click on Databases
 3. Create Database
 4. Choose a Database Creation Method - Standard create
 5. Select Latest Engine Version
 6. Choose Free Tier Template (for educational purposes)
 7. Provide Your Master Username (DB name)
 8. Credentials Management - Choose self-managed (create your own password) and store the credentials in AWS Secrets Manager.
 9. Public Access - Yes (strictly not recommended in production cases; since we are connecting to our RDS from outside the network, select yes. Ideally, in production, this should be 'NO').

For the rest of the settings, keep them as default.

The screenshot shows the Amazon RDS Databases page. The left sidebar includes links for Dashboard, Databases (which is selected and highlighted in blue), Query Editor, Performance insights, Snapshots, Exports in Amazon S3, Automated backups, Reserved instances, Proxies, Subnet groups, Parameter groups, Option groups, Custom engine versions, Zero-ETL integrations (with a 'New' badge), Events, and Event subscriptions.

The main content area displays a table titled "Databases (1)". The table has columns: DB identifier, Status, Role, Engine, Region & AZ, Size, and Recommendations. A search bar at the top says "Filter by databases". A "Create database" button is located in the top right corner. The table shows one database entry:

DB identifier	Status	Role	Engine	Region & AZ	Size	Recommendations
mysql-aws-de-db	Available	Instance	MySQL Community	ap-south-1a	db.t3.micro	

Keep note of the endpoint, port, username, and password - these are the things needed to establish a connection with the database.

Step 2: Store the RDS Credentials in AWS Secrets Manager

To avoid publicly exposing the credentials, securely store the RDS credentials in AWS Secrets Manager.

1. Select the AWS Secrets Manager service.
 2. Select Store a new secret.
 3. For Secret type, choose Credentials for Amazon RDS.
 4. Input your credentials.
 5. Choose the associated database.
 6. Give the secret a name.
 7. Click Store.

Secrets		
<input type="text"/> Filter secrets by name, description, tag key, tag value, owning service or primary Region		
Secret name	Description	Last retrieved (UTC)
dev-mysql-creds	-	-

Step3: Configuring Security Group for Database Access

We need to open the port on which the database service is running to allow inbound traffic. Our RDS is running in a VPC, which has a security group. Within this security group, we need to allow inbound traffic on port 3306, the port on which the database is running.

- Under the Connectivity & security section of the database instance, click on the VPC link.
- Click on the security group associated with the VPC.
- Click Edit inbound rules.
- Click Add rule.
- For Type, select MYSQL/Aurora.
- For Port range, enter 3306.
- For Source, enter 0.0.0.0/0 (to allow any IP).
- Click Save rules.

Inbound rules (2)					
<input type="text"/> Search					
	Name	Security group rule...	IP version	Type	Protocol
<input type="checkbox"/>	-	sgr-0ff8cf998481f7406	-	All traffic	All
<input type="checkbox"/>	-	sgr-00f3ee5fb0e951cd9	IPv4	MySQL/Aurora	TCP

Step4: Connect to the Database and create target db and table

A. Programmatically

In this step, we connect to the RDS instance and run some scripts. We need to install the MySQL connectors(`pip install mysql-connector-python`) to use functions and methods that will set up the connectivity, and Boto3 to connect to AWS Secrets Manager to retrieve the RDS credentials.

- Fetch secrets from the Secrets Manager in JSON format (e.g., `{ "username": "xxxxxx", "password": "xxxxx" }`).
- Establish a connection with the RDS using the MySQL connector.
- Execute our queries.

Check out the full code [[here](#)].

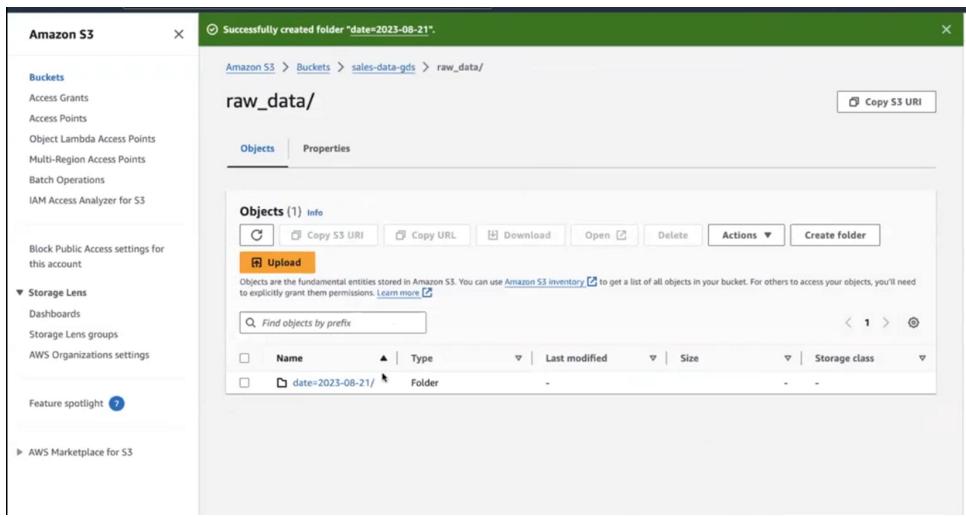
B. Through Terminal

You can also connect to your RDS instance using the terminal. Execute the following command to do so:
`mysql -h<host/endpoint> -P<port> -u<username> -p`

- Enter your password when prompted.
- After connecting, you can execute your SQL queries. [[Link to the queries](#)]

Step5: Setup the Data Lake - create an S3 Bucket for Raw Data

Create an S3 bucket where the client can upload CSV data in Hive partitions (using the date as the partition key).



Step 5: Create a Glue Catalog

Create a Glue catalog to serve as a central repository of metadata. Since we need metadata for both the source and target, we will create two tables in the catalog.

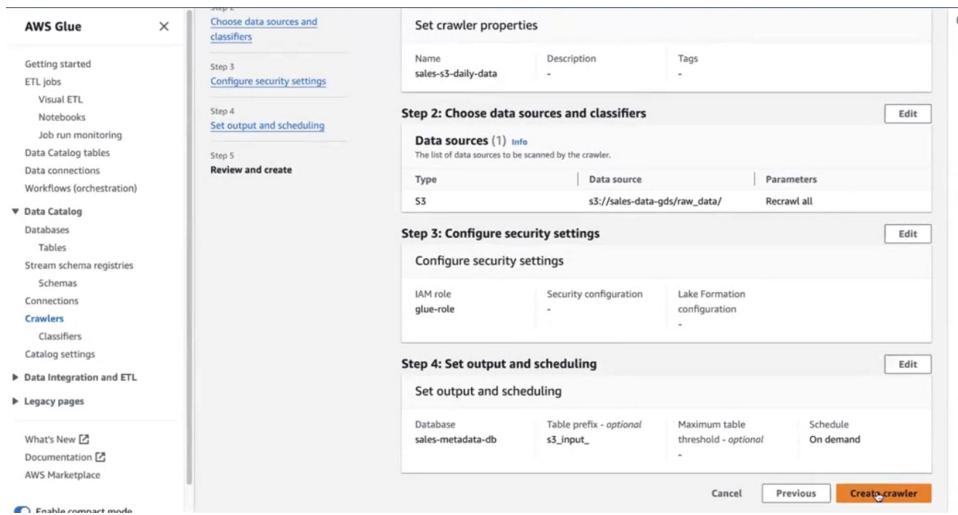
- Under the Data Catalog section in AWS Glue, select Databases.
- Give a name to the database.
- Click on create.

The screenshot shows the AWS Glue console. The left sidebar has sections for Getting started, ETL jobs, Data Catalog tables, Data connections, Workflows (orchestration), Data Catalog (with sub-options like Databases, Tables, Stream schema registries, etc.), and Data Integration and ETL. The main area is titled "AWS Glue > Databases". It shows a table for "Databases (1)". The table has one row with the database name "sales-metadata-db". The table includes columns for Name, Description, Location URI, and Created on (UTC).

Step 6: Create a Crawler for the Source

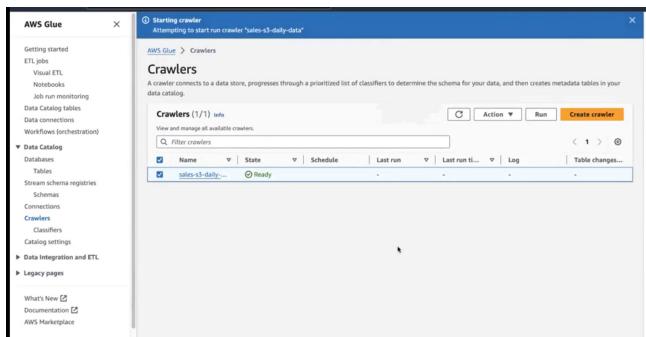
Create a crawler to fetch the metadata of the source/raw data in S3 and store it in a table in the Glue Catalog. For example, Crawler identifies the partition name, columns present in the partition, and their data types. Once we have this metadata, we can read data through the catalog instead of directly from S3 each time.

- Select Crawlers in AWS Glue.
- Click Create crawler.
- Give the crawler a name.
- For data source configuration, choose Is your data already mapped to Glue tables? - Not yet.
 - Click Add a data source.
 - Choose S3 as the source.
 - Select the source S3 bucket.
 - Select the path of the source folder (choose the top folder, not the file).
 - Select the Crawl all sub-folders option.
- For IAM Role, since this crawler scans S3, it needs permission to access S3. Create a new role for this glue and add the S3 read access policy.
- Select the previously created database in the Glue Catalog to store the source S3 metadata as a table.
- Give the table a name.
- Choose the On demand option for the crawler schedule. (We can automate this crawler as well, but for better understanding of things lets keep it ondemand)



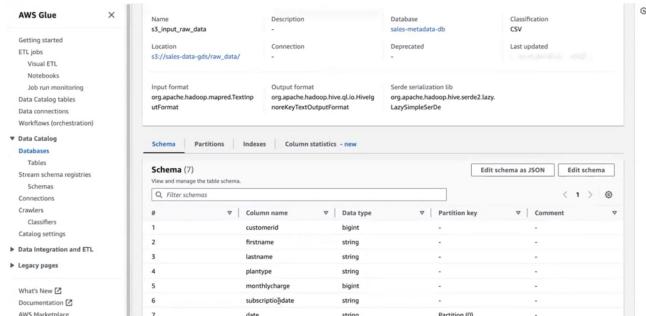
Step 7: Run the Crawler to Create the Schema

Once the crawler is created, select it and run it.



After the crawler runs successfully, follow these steps:

- Open the Databases under the Data Catalog.
- Select the database you created.
- You will see a new table reflected there. Click on it to view the schema that the crawler has created.
- If you want to modify the schema, use the Edit schema option.



Step 8: Create a JDBC connection to RDS.

Once our source metadata is set up in the catalog, we need to establish a similar metadata table for the target. This facilitates seamless integration within the Glue ETL pipeline. For our target, which is RDS, we must create a connection to it. This connection will be utilized in our crawler to access the table.

To create the connection:

- Navigate to "Connections" under Data Catalog.
- Choose "Create connection."
- Select "MySQL" as our RDS is MySQL.
- Choose the RDS instance for the connection.
- Specify the database you want to crawl in the RDS.
- Input RDS credentials.
- Provide a name and create the connection.

After creating the connection, it's essential to test it:

- Select the created connection.
- Under actions, click "Test connection."
- Grant permission to connect to RDS by selecting the previous IAM role and adding the RDS policy to it.
- Click "Confirm." If the connection is successful, you will receive confirmation.

The screenshot shows the AWS Glue Connections page. On the left, there's a sidebar with navigation links like 'Getting started', 'ETL jobs', 'Data Catalog tables', 'Data connections', 'Connections', 'Data integration and ETL', and 'Legacy pages'. The main area has two tabs: 'Marketplace connectors' and 'Custom connectors'. Under 'Marketplace connectors', there's a button 'Go to AWS Marketplace'. Under 'Custom connectors', there's a button 'Create custom connector'. Below these tabs is a table titled 'Connectors (0) Info' with columns 'Name', 'Type', and 'Last modified'. A message says 'No connectors' and 'No connectors to display. You can create a custom connector or get a Marketplace connector.' At the bottom is a 'Create custom connector' button. Another tab 'Connections (3) Info' is shown below, with a table listing three connections: 'Mysql New Connection' (selected), 'Mysql Connection', and 'Redshift-New-Connection'. The 'Mysql New Connection' row includes 'Actions' (View details, Delete, Edit), 'Type' (JDBC), and 'Last modified' (Mar 16, 2024).

The screenshot shows the 'Test Connection' dialog box. It displays a message 'Successfully connected to the data store with connection Mysql New Connection.' with a green checkmark icon. There are 'View log' and 'Cancel' buttons at the bottom. In the background, the AWS Glue Connections page is visible, showing the same list of connections as the previous screenshot.

Step 9: Endpoint creation

To ensure secure interaction between Glue, RDS, and S3, it's crucial to establish an endpoint. Glue relies on S3 for temporary data storage and metadata exchange over the network. This secure endpoint enables Glue to interact with S3, facilitating proper connection establishment.

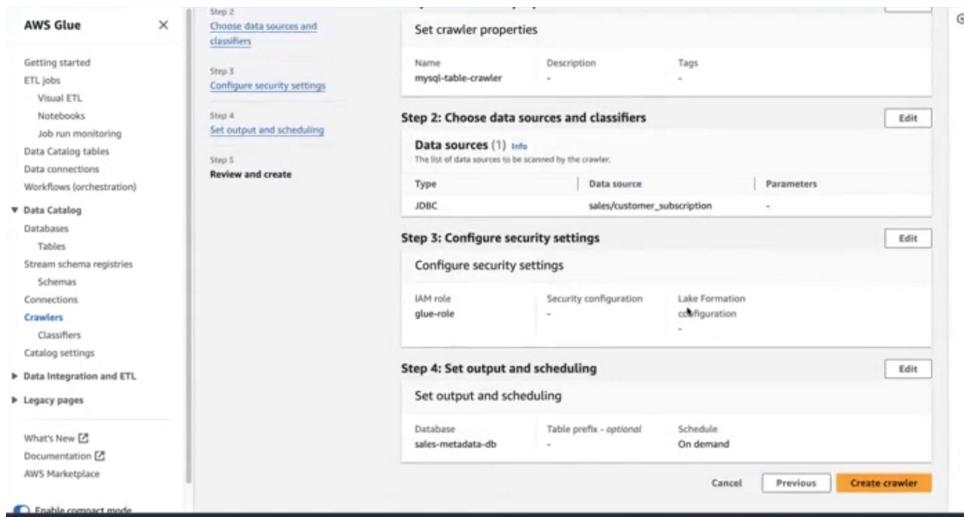
To create the endpoint:

- Navigate to VPC under the networking section in RDS.
- Click on "Endpoints."
- Choose "Create endpoint."
- Provide a name.
- Since we're creating an endpoint to S3, select S3 in the services and choose the one ending with ".s3" (not any other) and select the type as gateway.
- Select the same VPC where our RDS resides.
- Choose the all routing table.
- Apply the "Allow full access" policy.
- Click on "Create endpoint" to finalize the process.

Step10: Create crawler for target table

Now that we've established our JDBC connection and created an endpoint, it's time to set up a crawler to extract data from our target table in RDS.

- Start by creating a new crawler.
- Provide a name for the crawler.
- For the data source, select the JDBC connection we previously set up.
- Specify the table you wish to crawl, ensuring the format is "databasename tablename."
- Add the JDBC connection by clicking the corresponding button.
- Attach the Glue IAM role we created earlier.
- Choose the "On demand" option.
- Finally, initiate the creation of the crawler.



After creating the crawler, execute it to see the schema on the catalog, similar to what we did for the source crawler. Upon successful execution, the generated schema will be available.

The screenshot shows the AWS Glue Data Catalog interface. The left sidebar has the same navigation as the previous screenshot. The main area shows a table with the following schema:

#	Column name	Data type	Partition key	Comment
1	monthlycharge	int	-	-
2	subscriptiondate	date	-	-
3	firstname	string	-	-
4	plantype	string	-	-
5	customerid	int	-	-
6	lastname	string	-	-

Step 11: Create an ETL Job Using Visual ETL

To handle incremental reads, we'll use Job Bookmarking. This feature keeps track of all the files read so far and looks for newly created objects.

ETL consists of three parts: Extract, Transform, and Load.

Part 1: Extract Data from Source

- After selecting the Visual ETL option, provide a name for the job.
- Click on 'Add Nodes.'
- Under Source, choose Glue Data Catalog (this is where our metadata is stored).
- In the Data Source Properties tab, provide a name for the data source. Choose the table and database from which the metadata will be read, and attach the Glue IAM role.
- The system will automatically start fetching the data for you.

The screenshot shows the AWS Glue Visual ETL interface. The top bar has tabs for 'Visual', 'Script', 'Job details', 'Runs', 'Data quality - updated', 'Schedules', and 'Version Control'. A modal window titled 'Data source properties - Data Catalog' is open, showing the following configuration:

- Name: S3_Data_Source
- Database: sales-metadata-db
- Table: s3_input_raw_data
- Use runtime parameters: checked

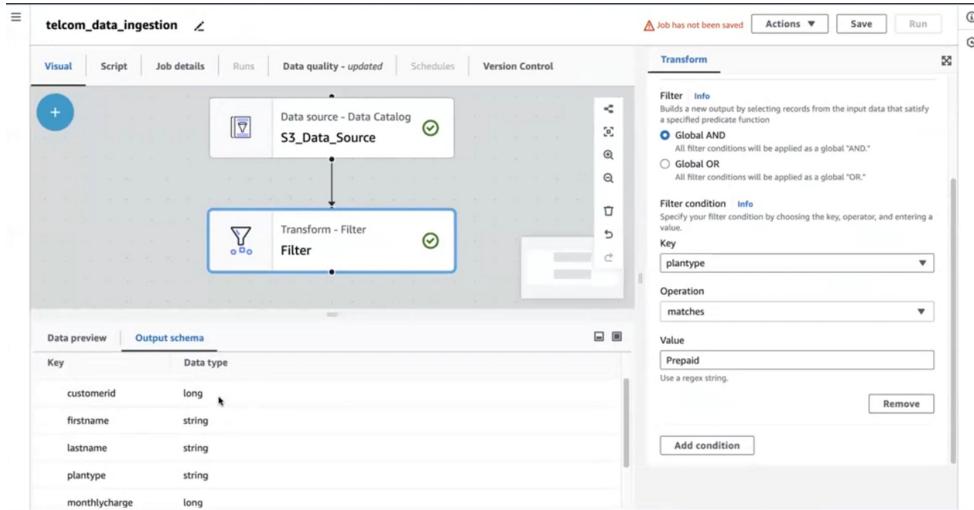
Below the modal, the 'Schema' section shows the extracted schema:

Key	Data type
customerid	long
firstname	string
lastname	string
plantype	string
monthlycharge	long
subscriptiondate	string

Part2: Transform The Data

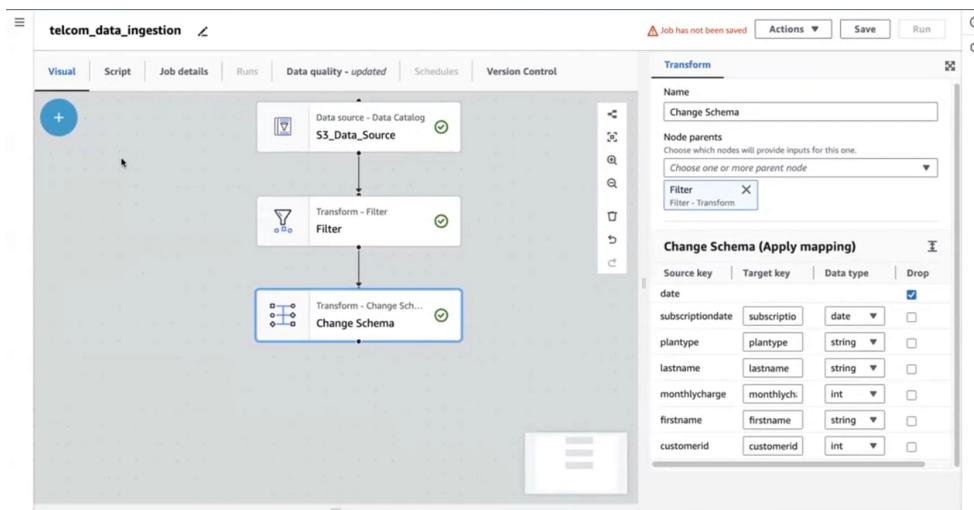
- Click on 'Add Node' under the Transform section and choose the Filter option.
- Provide a name for the transformation and add a filter condition:

- c. a. Filter condition: Key: Plantype, Operation: matches, Value: Prepaid
d. Once saved, the filtered data schema will be displayed.
e. After the output schema is generated for the filter stage, review the properties and make any necessary edits if there are any mismatches.



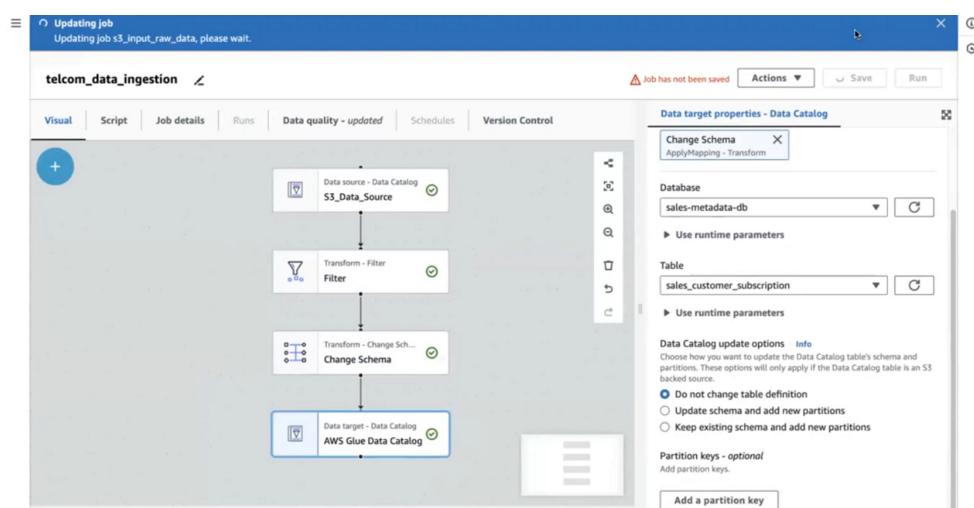
In our case, customerId is read as long instead of int, monthlyCharge as long instead of int, and subscriptionDate as string instead of date. Additionally, an extra column 'date' is present. We can use the change schema property under the Transform section to rectify this.

After selecting, you can change the key name, data type, or drop the column altogether as per requirements. Remember, the final schema of the Transform section should match the schema of the target table.



Part3: Add Target

- Choose "Data Catalog" as our target is defined in the catalog.
- Select the database and target table.
- Click the "Save" button.

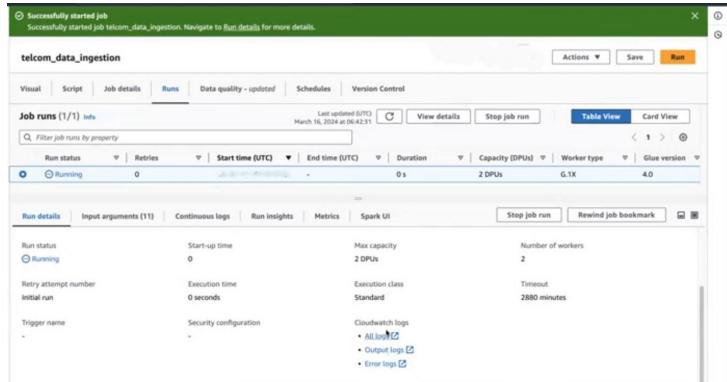


Step12: Define the Glue Infrastructure for Computation

Since Glue performs distributed computation, we need to select the runtime infrastructure. When the ETL job is executed, it will use this infrastructure to perform the computation.

- a. Select the "Job details" tab.
- b. Choose the Spark 3.3 version.
- c. You can either choose the number of worker nodes or use the auto-scaling feature.
- d. Enable job bookmarking for incremental reads.
- e. Set the number of retries if the ETL job fails at any point, and configure the timeout settings.
- f. Under the advanced settings, find the path to the script. This script is distributed to the nodes for computation.
- g. Now, run your ETL pipeline by clicking on "Run."
- h. Check the "Runs" section to see the execution process and the parameters passed.

This ensures that your Glue job is properly configured for efficient and reliable execution.



Once completed, query the database and check the results.

```
mysql> select * from customer_subscription;
ERROR 2013 (HY000): Lost connection to MySQL server during query
No connection. Trying to reconnect...
Connection id: 58
Current database: sales

+-----+-----+-----+-----+-----+
| CustomerID | FirstName | LastName | PlanType | MonthlyCharge | SubscriptionDate |
+-----+-----+-----+-----+-----+
| 1102 | Abigail | Johnson | Prepaid | 30 | 2022-02-15 |
| 1104 | Addison | Jones | Prepaid | 28 | 2022-05-05 |
| 1106 | Aiden | Davis | Prepaid | 35 | 2022-08-22 |
| 1108 | Alex | Miller | Prepaid | 32 | 2023-02-28 |
| 1110 | Alexander | Rodriguez | Prepaid | 29 | 2023-05-21 |
| 1112 | Alicia | Anderson | Prepaid | 31 | 2023-07-15 |
| 1114 | Alvin | Martinez | Prepaid | 30 | 2023-08-15 |
| 1116 | Amber | Lee | Prepaid | 28 | 2023-10-05 |
| 1118 | Amy | Gonzalez | Prepaid | 32 | 2023-11-11 |
| 1120 | Andrea | Lewis | Prepaid | 33 | 2023-12-24 |
+-----+-----+-----+-----+-----+
10 rows in set (2.23 sec)

mysql>
```

The 2023-08-21 CSV file contained 10 prepaid customers, which is now reflected in the database. If you want to check new folders or files, make sure to load the data onto S3, run the crawlers, and then execute the ETL pipeline, as we have configured them as on-demand functions.