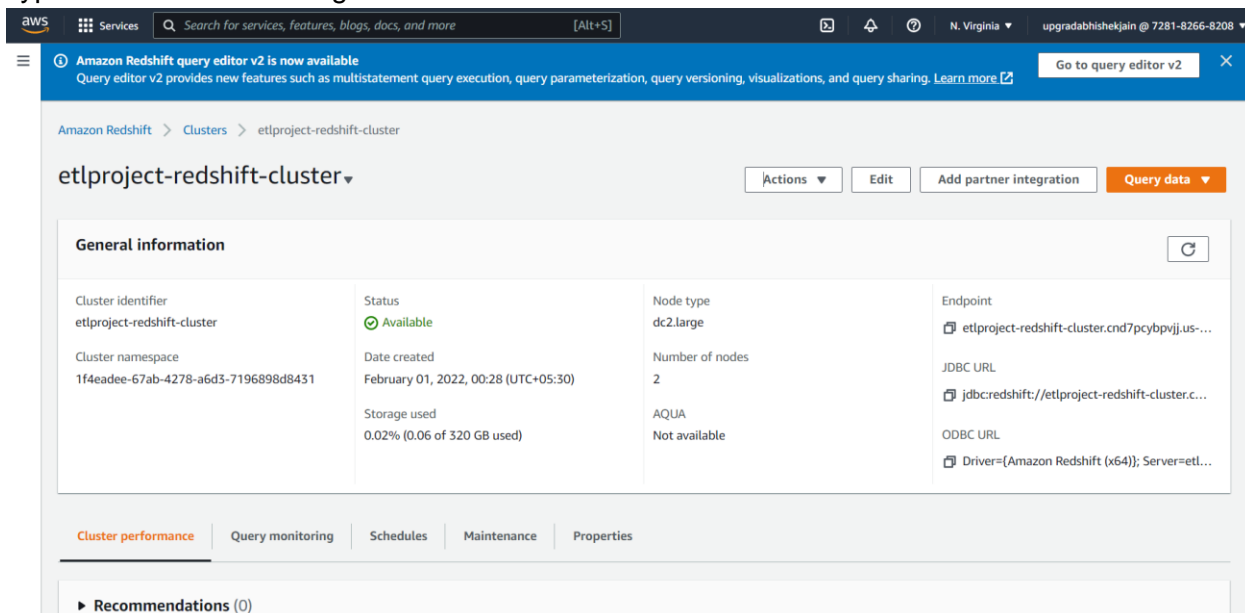# Creation of a Redshift Cluster

**Screenshots of the configuration of the Redshift cluster that you have created:**
Type of machine : dc2.large
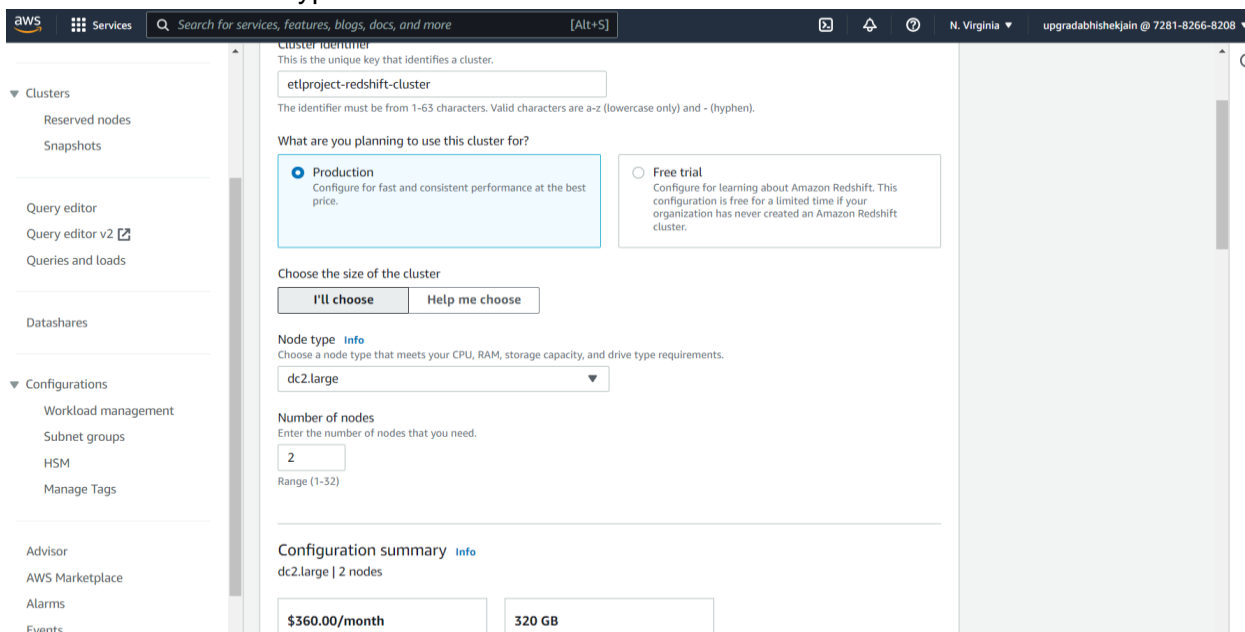


Screenshot of various configurations while setting up Red shift cluster
Cluster use and node type

# Database Configurations



# IAM Roles

Network and security: selecting VPC and subnet group



Configuring DB port

# Setting up a database in the Redshift cluster and running queries to create the dimension and fact tables

**View all the data in Amazon S3 Bucket:**



**Queries to create the various dimension and fact tables with appropriate primary and foreign keys:**

**Create a schema for dimension and fact table**

- **create schema atm_data;**

- **Create Dim_location table**
  create table atm_data.DIM_LOCATION
  (

          location_id int not null DISTKEY SORTKEY,
          location varchar(50),
          streetname varchar(255),
          street_number int,
          zipcode int,
          lat decimal(10,3),
          lon decimal(10,3),
          PRIMARY KEY(location_id)
  );



- **Create Dim_ATM table**
  create table atm_data.DIM_ATM
  (

          atm_id int not null DISTKEY SORTKEY,
          atm_number varchar(20),
          atm_manufacturer varchar(50),
          atm_location_id int,
          PRIMARY KEY(atm_id),
          FOREIGN KEY(atm_location_id) references
  atm_data.DIM_LOCATION(location_id)
  );

- **Create Dimension date table**

```
create table atm_data.DIM_DATE
(
        date_id int not null DISTKEY SORTKEY,
        full_date_time timestamp,
        year int,
        month varchar(20),
        day int,
        hour int,
        weekday varchar(20),
        PRIMARY KEY(date_id)
);
```

- **Create card type dimension table**
  ```
  create table atm_data.DIM_CARD_TYPE
  (
          card_type_id int not null DISTKEY SORTKEY,
          card_type varchar(30),
          PRIMARY KEY(card_type_id)
  );
  ```

- **Create atm transaction fact table**
  ```
  create table atm_data.FACT_ATM_TRANS
  (
          trans_id bigint not null DISTKEY SORTKEY,
          atm_id int,
          weather_loc_id int,
          date_id int,
          card_type_id int,
          atm_status varchar(20),
          currency varchar(10),
          service varchar(20),
          transaction_amount int,
          message_code varchar(225),
          message_text varchar(225),
          rain_3h decimal(10,3),
          clouds_all int,
          weather_id int,
          weather_main varchar(50),
          weather_description varchar(225),
          PRIMARY KEY(trans_id),
          FOREIGN KEY(weather_loc_id) references
  atm_data.DIM_LOCATION(location_id),
          FOREIGN KEY(atm_id) references atm_data.DIM_ATM(atm_id),
          FOREIGN KEY(date_id) references atm_data.DIM_DATE(date_id),
          FOREIGN KEY(card_type_id) references
  atm_data.DIM_CARD_TYPE(card_type_id)
  );
  ```
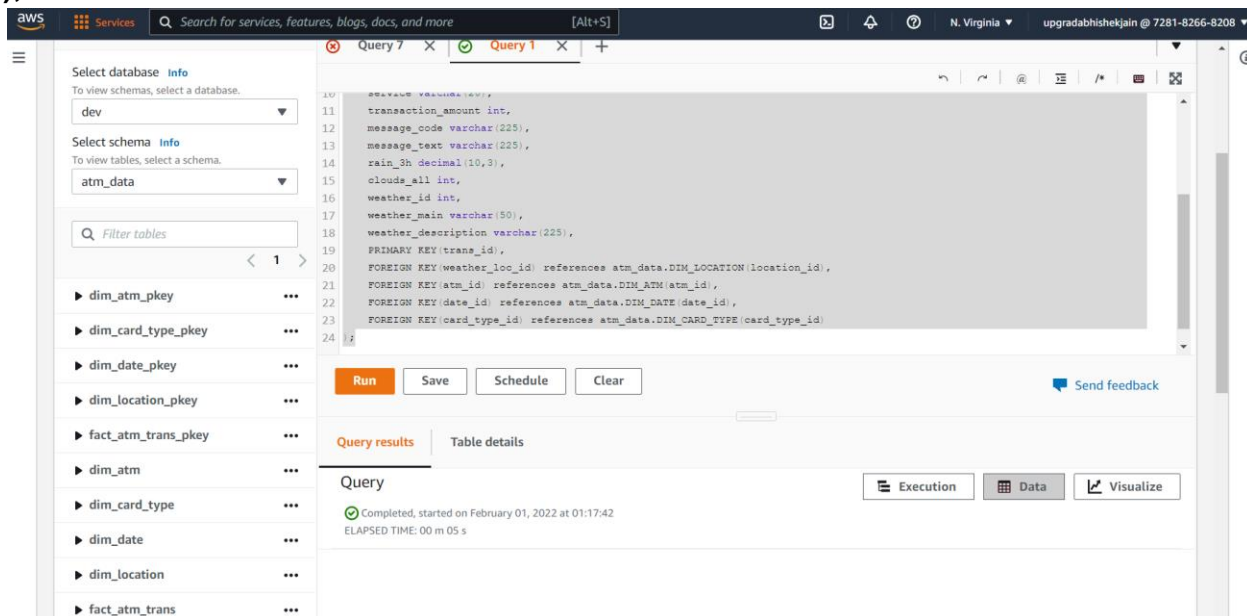
# Loading data into a Redshift cluster from Amazon S3 bucket

**Queries to copy the data from S3 buckets to the Redshift cluster in the appropriate tables**
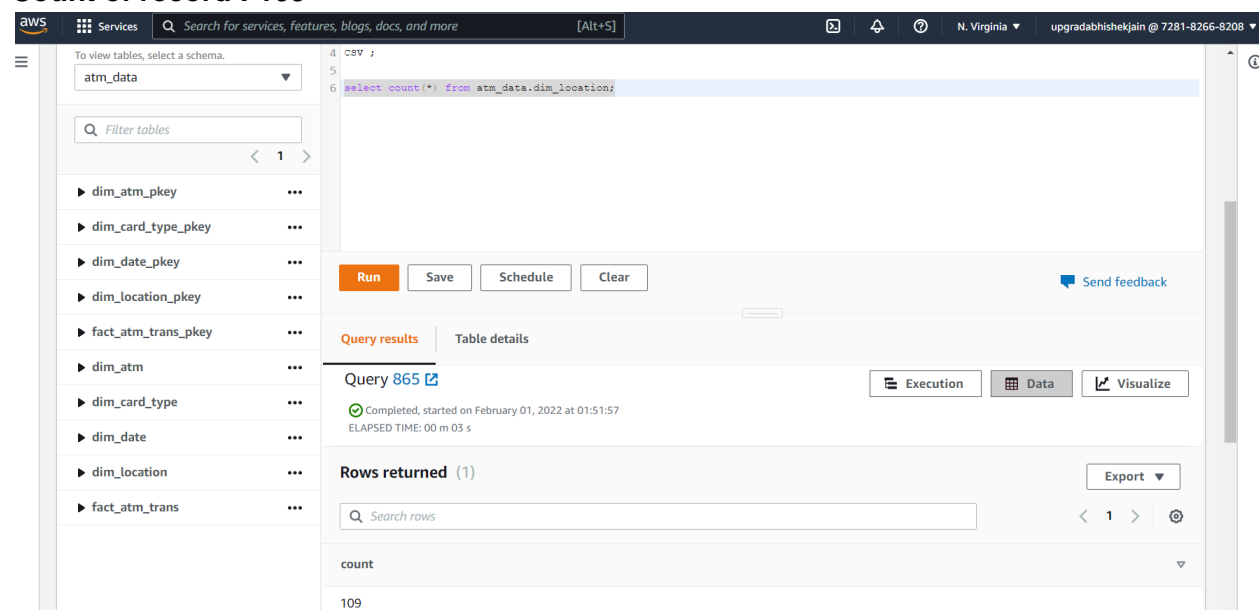
- **Copy data to DIM_LOCATION and DIM_ATM table**

copy atm_data.dim_location from 's3://etlprojectbyabhishek/dim_location/part-00000-2c3f4932-5d7f-4a0f-b030-9b3283149f77-c000.csv'
iam_role 'arn:aws:iam::728182668208:role/upgrad_redshift_s3'
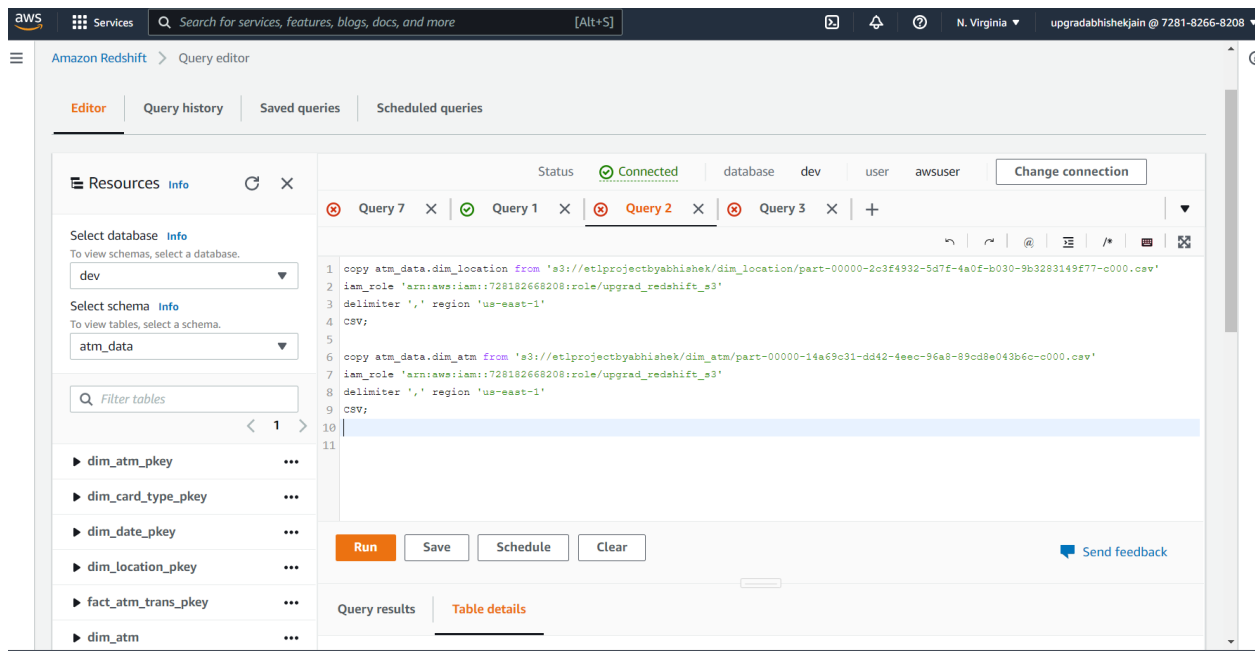delimiter ',' region 'us-east-1'
CSV;

**Count of record : 109**



copy atm_data.dim_atm from 's3://etlprojectbyabhishek/dim_atm/part-00000-14a69c31-dd42-4eec-96a8-89cd8e043b6c-c000.csv'
iam_role 'arn:aws:iam::728182668208:role/upgrad_redshift_s3'
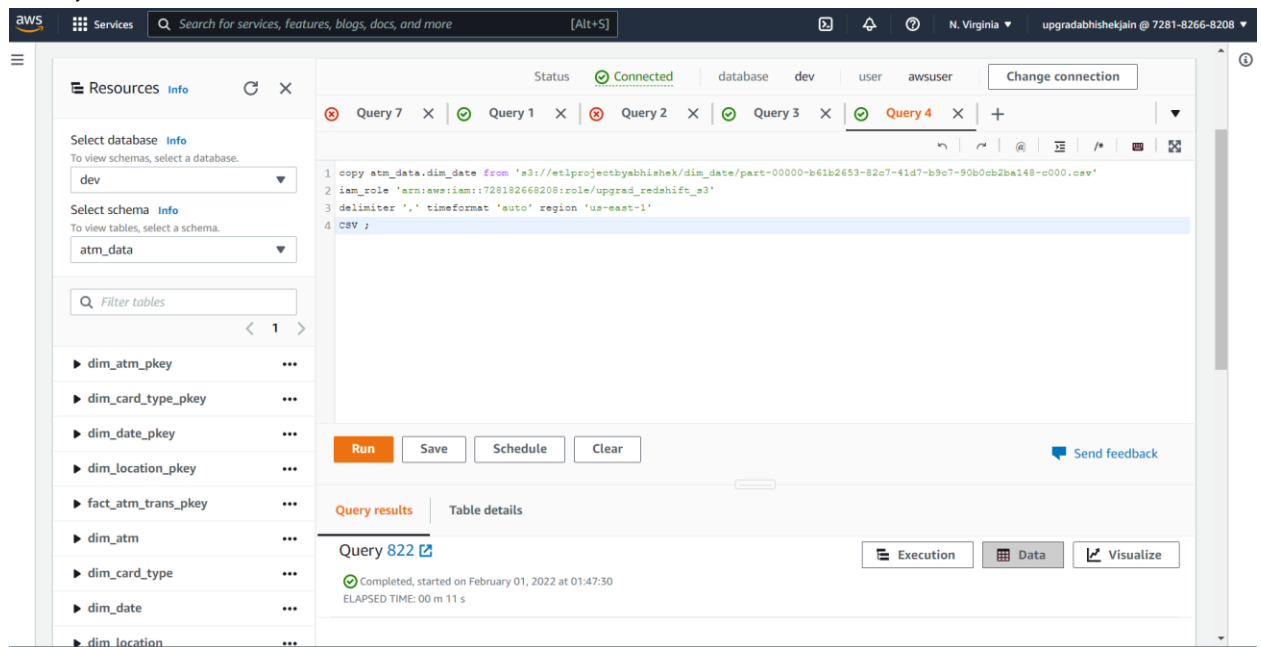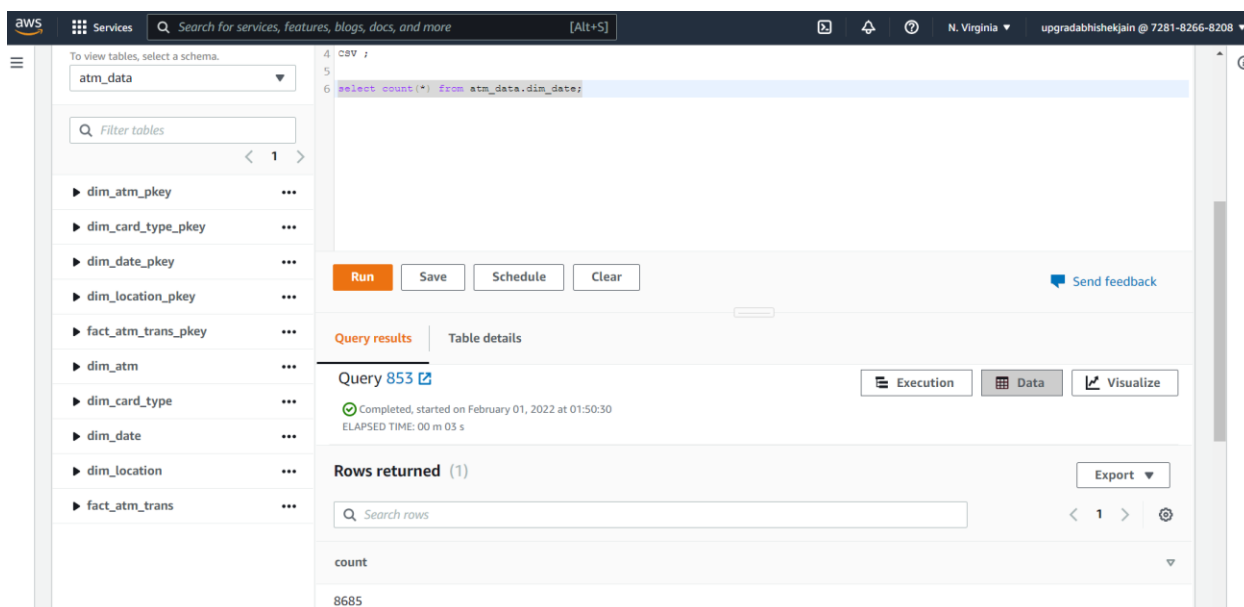delimiter ',' region 'us-east-1'
CSV;

- **Copy data to DIM_DATE table**
  copy atm_data.dim_date from 's3://etlprojectbyabhishek/dim_date/part-00000-b61b2653-82c7-41d7-b9c7-90b0cb2ba148-c000.csv'
  iam_role 'arn:aws:iam::728182668208:role/upgrad_redshift_s3'
  delimiter ',' timeformat 'auto' region 'us-east-1'
  CSV ;



**Count of record :8685**

- **Copy data to DIM_CARD_TYPE table**
  **copy atm_data.dim_card_type from 's3://etlprojectbyabhishek/dim_card_type/part-00000-ca80c3fe-5cf2-433c-b05a-10fd2ead1952-c000.csv'**
  **iam_role 'arn:aws:iam::728182668208:role/upgrad_redshift_s3'**
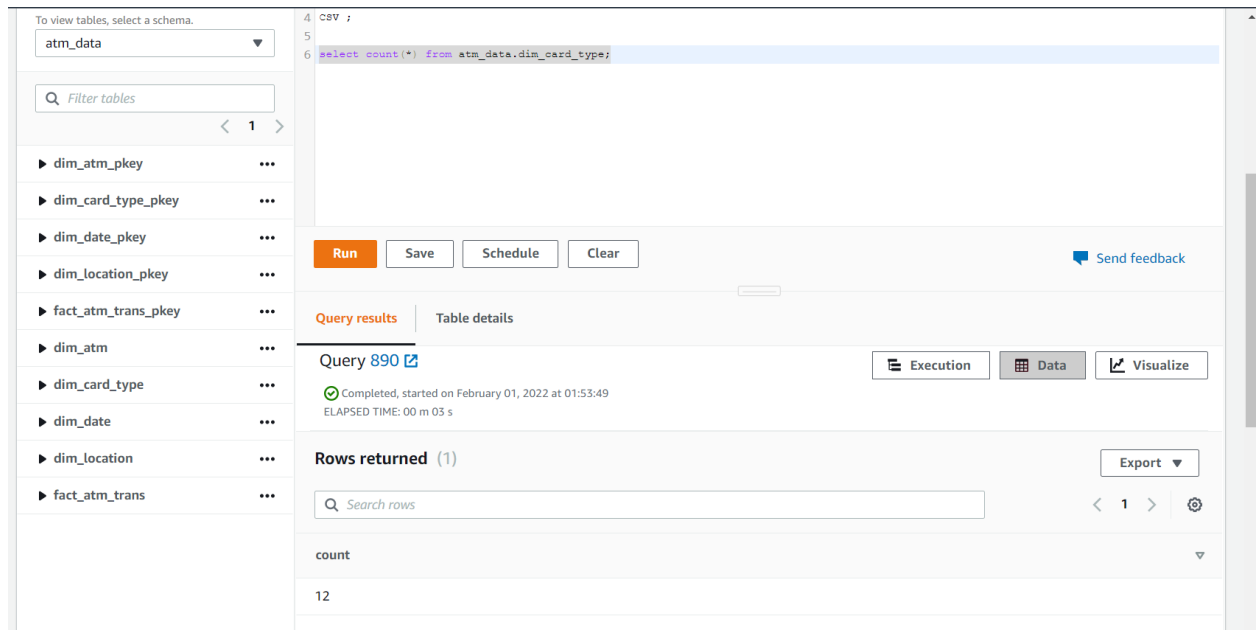  **delimiter ',' region 'us-east-1'**
  **CSV;**



**Count of record :12**
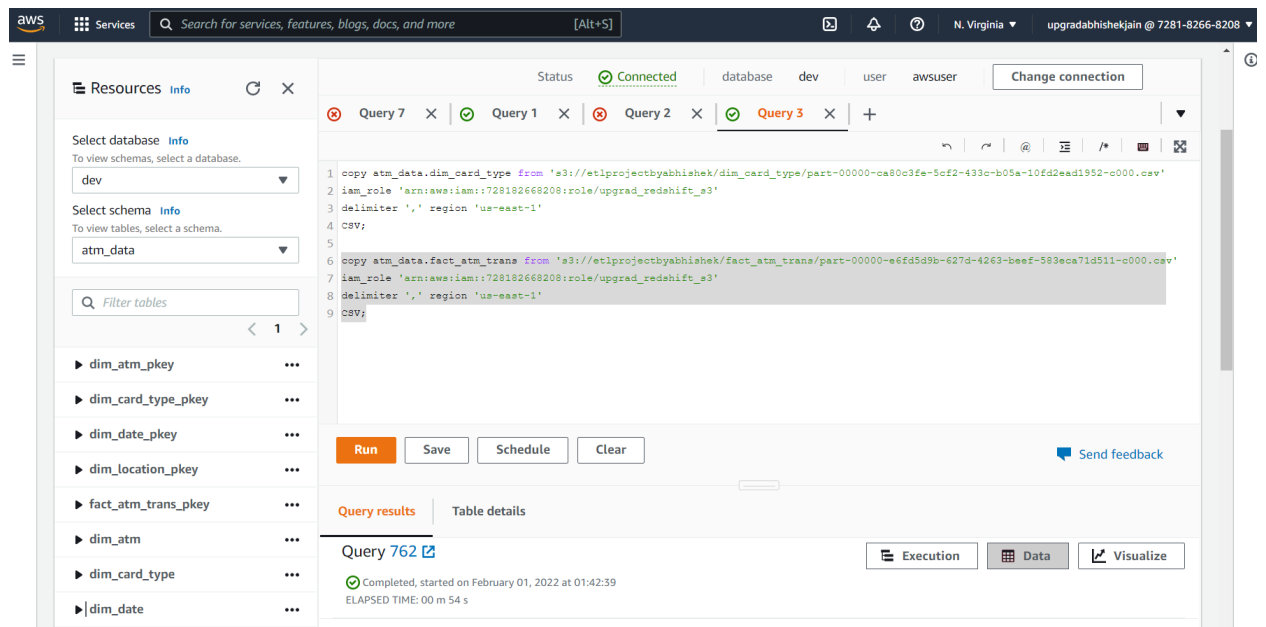
- **Copy data to FACT_ATM_TRANSACTION table**
  copy atm_data.fact_atm_trans from
  's3://etlprojectbyabhishek/fact_atm_trans/part-00000-e6fd5d9b-627d-4263-beef-583eca71d511-c000.csv'
  iam_role 'arn:aws:iam::728182668208:role/upgrad_redshift_s3'
  delimiter ',' region 'us-east-1'
  CSV;



**Count of record : 2468572**