## SEP 721 – Data Analytics, Machine Learning and AI on Cloud Platforms
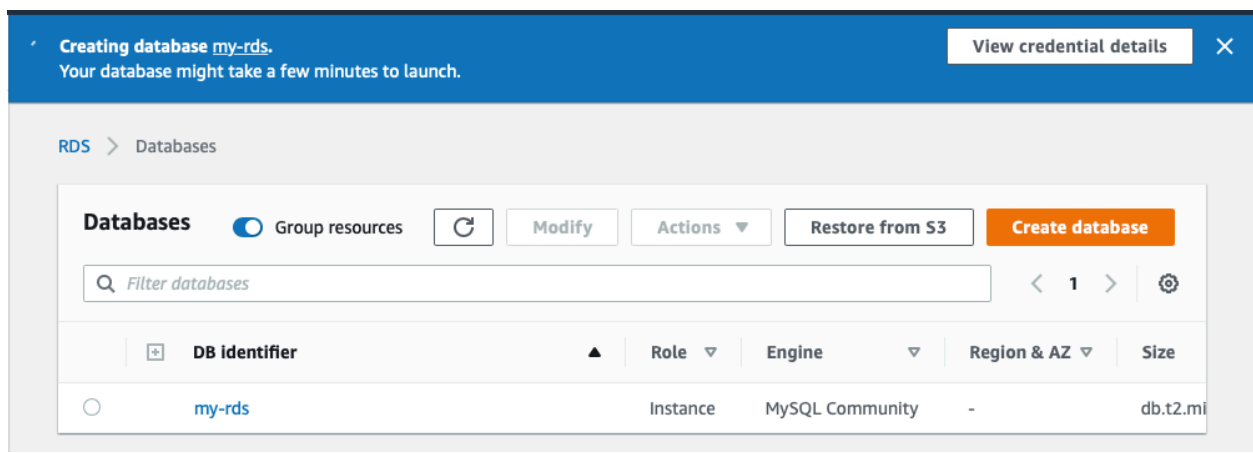
## Assignment 2: Qwiklabs- 1,2 and 3

Submitted by,
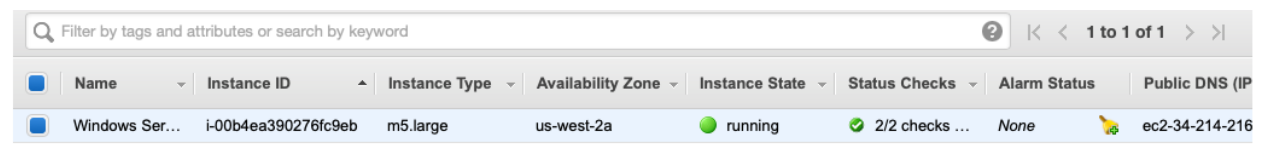
Greeshma Gopal(gopalg)

ID- 400245291

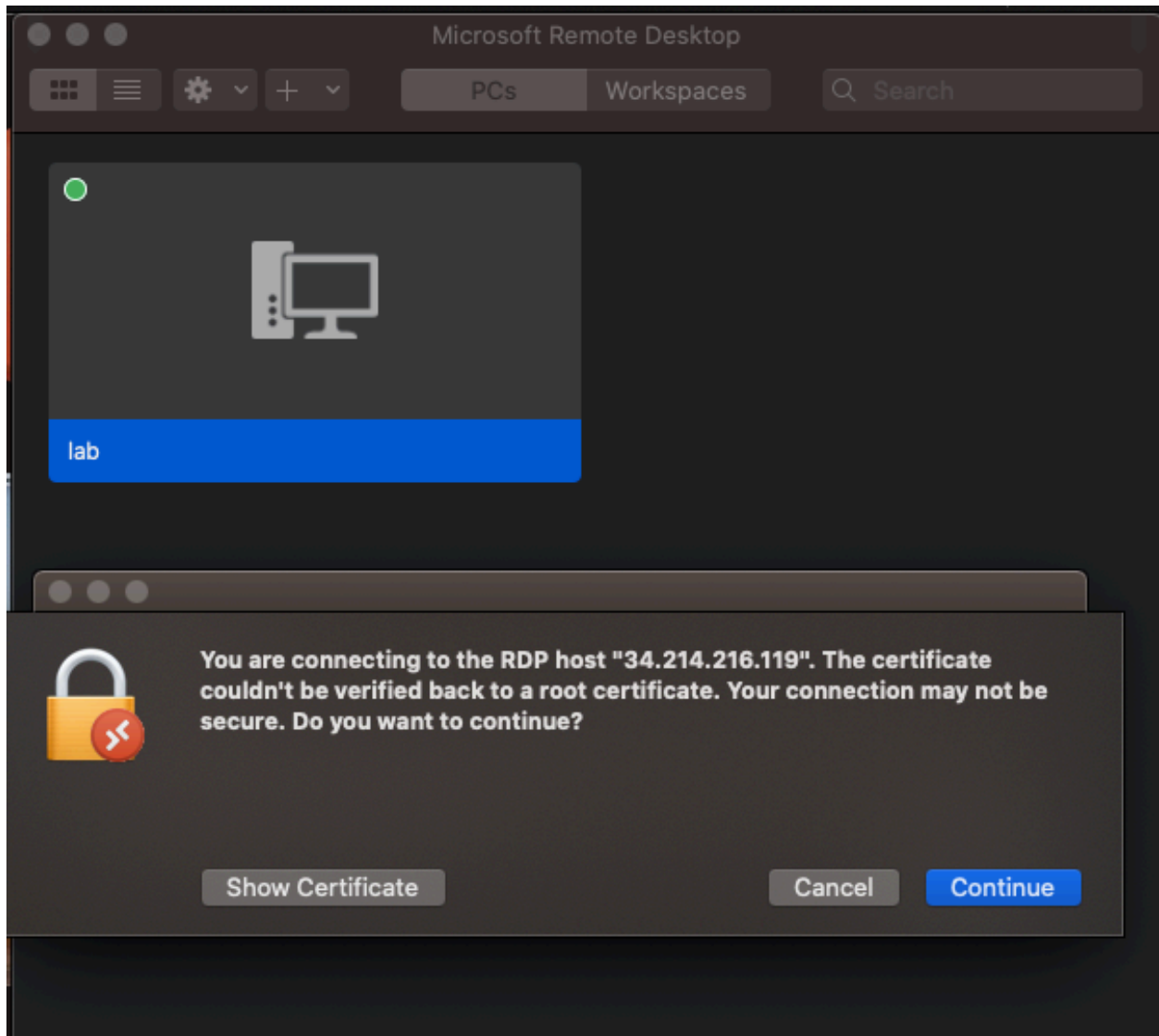# Lab 1: Introduction to Amazon Relational Database Service (RDS) (Windows)

- I have used qwiklab credits to get the login credentials. The database was created with the configuration provided

**Creating database my-rds.**
Your database might take a few minutes to launch.

View credential details ✕

RDS > Databases

**Databases**   ⬤ Group resources   ↻   Modify   Actions ▾   Restore from S3   **Create database**

Q Filter databases   ‹ 1 ›   ⚙

| ⊞ | DB identifier ▲ | Role ▽ | Engine ▽ | Region & AZ ▽ | Size |
|---|---|---|---|---|---|
| ○ | my-rds | Instance | MySQL Community | - | db.t2.mi |

- Logging in to EC2 instance. After the database was created, the EC2 instance which was created has to be logged in to.

Q Filter by tags and attributes or search by keyword   ❓ |‹ ‹ 1 to 1 of 1 › ›|

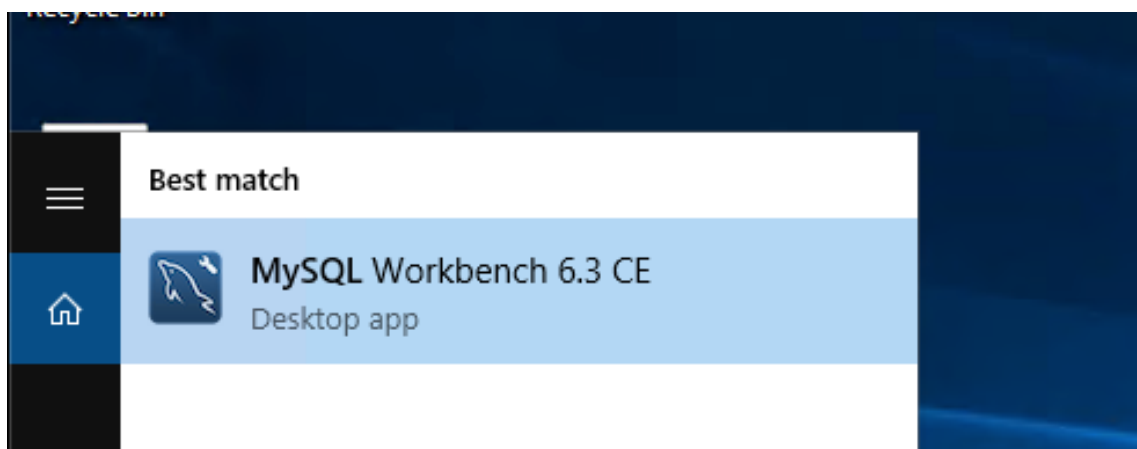| | Name | Instance ID ▲ | Instance Type ▾ | Availability Zone ▾ | Instance State ▾ | Status Checks ▾ | Alarm Status | Public DNS (IP |
|---|---|---|---|---|---|---|---|---|
| ☑ | Windows Ser... | i-00b4ea390276fc9eb | m5.large | us-west-2a | ⬤ running | ✔ 2/2 checks ... | None | ec2-34-214-216 |

- Post this, I downloaded the Microsoft remote desktop, this is to access the MySQL database.



- The desktop was successfully logged in.

- We will have to access MySQL in the desktop

- A connection has to be setup with the credentials provided

- The MySQL can be queried and utilized like a regular database

# Lab 2: Introduction to Amazon DynamoDB

- Create a new table in DynamoDB with a primary key and sort key for uniquely identifying each item

| Name | |
|------|---|
| ● Music | ▲ |

No CloudWatch alarms have been triggered for this table.

**Stream details**

| | |
|---|---|
| Stream enabled | No |
| View type | - |
| Latest stream ARN | - |

**Manage Stream**

**Table details**

| | |
|---|---|
| Table name | Music |
| Primary partition key | Artist (String) |
| Primary sort key | Song (String) |
| Point-in-time recovery | DISABLED  Enable |
| Encryption Type | DEFAULT  Manage Encryption |
| KMS Master Key ARN | Not Applicable |
| Encryption Status | |
| CloudWatch Contributor Insights | DISABLED  Manage Contributor Insights  PREVIEW |
| Time to live attribute | DISABLED  Manage TTL |
| Table status | Active |

- Adding data to the table



**Edit item**

Tree ▾

▼ Item {2}

  ● Artist String : Pink Floyd

  ●  ⊕ Insert ▾   g : Money

       String
       Binary
       Number
       StringSet
       NumberSet
       BinarySet
       Map
       List
       Boolean
       Null

  ⊗ Remove

- Editing to add more fields

**Edit item**

| Tree ▾ | ⇅ ⇊ |
|---|---|

▼ Item {3}
- ⊕    Album String : The Dark side of the Moon
- ⊕    Artist String : Pink Floyd
- ⊕    Song String : Money

- Creating multiple items

| | Artist ⓘ | ▲ | Song | ⌄ | Album | ⌄ | Year | ⌄ | Genre | ⌄ | LengthSeconds |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | John Lennon | | Imagine | | Imagine | | 1971 | | Soft rock | | |
| ☐ | Pink Floyd | | Money | | The Dark sid... | | 1973 | | | | |
| ☐ | Psy | | Gagnam style | | Psy 6 (Six Ru... | | 2011 | | | | 219 |

- Modifying existing item

**Edit item**

```
Tree ▾    ⬍  ⬍

      ▼ Item {5}
  ⊕       Album String : Psy 6 (Six Rules), Part 1
  ⊕       Artist String : Psy
  ⊕       LengthSeconds Number : 219
  ⊕       Song String : Gagnam style
  ⊕       Year Number : 2012
```

- Querying the table

**Scan: [Table] Music: Artist, Song** ^                    Viewing 1 to 3 items

| Query ⬍ | [Table] Music: Artist, Song | ⬍ | ^ |

| **Partition key** | Artist | String | = | Psy |
| **Sort key** | Song | String | = ⬍ | Gangnam Style |
| **Filter** | Enter attribute | String ⬍ | = ⬍ | Enter value | ✕ |
| | ⊕ Add filter | | | |

**Sort**  ⦿ Ascending  ◯ Descending
**Attributes**  ⦿ All  ◯ Projected

| | Artist ⓘ | ▲ | Song | ▾ | Album | ▾ | Year | ▾ | Genre | ▾ | LengthSeconds |

- Scanning the table

- Deleting the table

## Lab 3 : Working with Amazon Redshift

- Creating clusters using amazon redshift

## Cluster details

**Cluster identifier**
This is the unique key that identifies a cluster.

lab

The identifier must be from 1-63 characters. Valid characters are a-z (lowercase only) and - (hyphen).

**Database port (optional)**
Port number of the port where the database accepts inbound connections.

5439

The port must be numeric (1150-65535).

**Master user name**
Enter a login ID for the master user of your DB instance.

master

The name must be 1-128 alphanumeric characters, and it can't be a **reserved word**.

**Master user password**

••••••••••

☐ Show password

## ▼ Network and security

**Virtual private cloud (VPC)**
This VPC defines the virtual networking environment for this cluster.

Lab VPC
vpc-04d4845d28016ddef ▼

ⓘ You can't change the VPC associated with this cluster after the cluster has been created. Learn more ↗     ✕

**VPC security groups**
This VPC security groups define which subnets and IP ranges the cluster can use in the VPC.

Choose one or more security groups ▼

Redshift Security Group ✕
sg-09eb8761a0ffe4e0d

**Cluster subnet group**
Choose the Amazon Redshift subnet group to launch the cluster in.

qls-12025451-0f6a6502ed005f6e-redshiftclustersubnetgroup-1lmqhvl5g760i ▼

Availability Zone

- Cluster has been created with the above configuration

- Connecting to the cluster using pgweb IP which was provided



- The table flight was created post which the number of rows in the table is being fetched

- The data has been copied to the new table which was created. The below query will fetch the data from the table

- In the below query we are grouping the result based on the carrier. The top three carriers based on number of departures are Southwest Airlines Co., Delta Airlines Co., and American Airlines Inc.

```
1   SELECT
2     carrier,
3     SUM (departures)
4   FROM flights
5   GROUP BY carrier
6   ORDER BY 2 DESC
7   LIMIT 10;
```

Run Query    Explain Query

| carrier | sum |
| --- | --- |
| Southwest Airlines Co. | 19846786 |
| Delta Air Lines Inc. | 17480331 |
| American Airlines Inc. | 15468460 |
| United Air Lines Inc. | 13402737 |
| Northwest Airlines Inc. | 9522027 |
| Continental Air Lines Inc. | 7929064 |
| US Airways Inc. | 7865177 |
| American Eagle Airlines Inc. | 6687299 |
| ExpressJet Airlines Inc. | 6270410 |
| USAir | 6052546 |

- The top three carriers based on number of passengers are Delta Airlines Co., Southwest Airlines Co., and American Airlines Inc.

```
1  SELECT
2    carrier,
3    SUM (passengers)
4  FROM flights
5  GROUP BY carrier
6  ORDER BY 2 DESC
7  LIMIT 10;
```

[Run Query]  [Explain Query]

| carrier | sum |
|---|---|
| Delta Air Lines Inc. | 1894041955 |
| Southwest Airlines Co. | 1812013508 |
| American Airlines Inc. | 1599088928 |
| United Air Lines Inc. | 1396240286 |
| Northwest Airlines Inc. | 846994753 |
| Continental Air Lines Inc. | 751752616 |
| US Airways Inc. | 738513678 |
| USAir | 406251028 |
| America West Airlines Inc. | 327787337 |
| Alaska Airlines Inc. | 290615980 |

- The top three carriers based on miles flown are United airlines Inc., Southwest Airlines Co., and American Airlines Inc.

| Rows | Structure | Indexes | Constraints | SQL Query | History | Activity | Connection |
|---|---|---|---|---|---|---|---|

```
1  SELECT
2    carrier,
3    SUM (miles)
4  FROM flights
5  GROUP BY carrier
6  ORDER BY 2 DESC
7  LIMIT 10;
```

[Run Query]  [Explain Query]

| carrier | sum |
|---|---|
| United Air Lines Inc. | 7478582456 |
| Delta Air Lines Inc. | 7356255471 |
| American Airlines Inc. | 6033117399 |
| Southwest Airlines Co. | 5326561689 |
| Continental Air Lines Inc. | 4454443574 |
| Northwest Airlines Inc. | 4277614432 |
| US Airways Inc. | 3193773517 |
| USAir | 2178906790 |
| America West Airlines Inc. | 1713042483 |
| Federal Express Corporation | 1514604035 |

- The top three carriers based on passenger-miles flown are Southwest Airlines Co., Delta Airlines Co., and American Airlines Inc.

```
1  SELECT
2    carrier,
3    SUM (passengers * miles)
4  FROM flights
5  GROUP BY carrier
6  ORDER BY 2 DESC
7  LIMIT 10;
```

Run Query    Explain Query

| carrier | sum |
| --- | --- |
| American Airlines Inc. | 1733245030322 |
| Delta Air Lines Inc. | 1618918706400 |
| United Air Lines Inc. | 1537952630588 |
| Southwest Airlines Co. | 1045280304865 |
| Continental Air Lines Inc. | 805455729207 |
| Northwest Airlines Inc. | 742067994834 |
| US Airways Inc. | 585542298085 |
| America West Airlines Inc. | 293772304506 |
| Alaska Airlines Inc. | 266206170484 |
| USAir | 240098568130 |

- The top three carriers based on freight pounds are Federal Express Corporation, United Parcel Service and Delta Airlines Inc.

```
1  SELECT
2    carrier,
3    SUM (freight_pounds)
4  FROM flights
5  GROUP BY carrier
6  ORDER BY 2 DESC
7  LIMIT 10;
```

Run Query    Explain Query

| carrier | sum |
| --- | --- |
| Federal Express Corporation | 115908103479 |
| United Parcel Service | 70144943595 |
| Delta Air Lines Inc. | 12727490652 |
| United Air Lines Inc. | 11310081837 |
| Korean Air Lines Co. Ltd. | 11237313800 |
| American Airlines Inc. | 10391736851 |
| China Airlines Ltd. | 9654006204 |
| Eva Airways Corporation | 6378076357 |
| Northwest Airlines Inc. | 6310003281 |
| Southwest Airlines Co. | 5494931618 |

- Creating a new table aircraft for joining

```
public
  Tables (2)
    aircraft
    flights
```

```
1  CREATE TABLE aircraft (
2    aircraft_code CHAR(3) SORTKEY,
3    aircraft       VARCHAR(100)
4  );
```

- Copying data from the bucket to the aircraft table

```
public
  Tables (2)
    aircraft
    flights
```

```
1  COPY aircraft
2  FROM 's3://us-west-2-aws-training/awsu-spl/spl-17/4.2.5.prod/data/lookup_aircraft.csv'
3  IAM_ROLE 'arn:aws:iam::723548672068:role/Redshift-Role'
4  IGNOREHEADER 1
5  DELIMITER ','
6  REMOVEQUOTES
7  TRUNCATECOLUMNS
8  REGION 'us-west-2';
```

- Selecting the data from the newly created table aircraft

```
1  SELECT *
2  FROM aircraft
3  ORDER BY random()
4  LIMIT 10;
```

Run Query    Explain Query

| aircraft_code | aircraft |
|---|---|
| 314 | Bell B-47J2 |
| 039 | Cessna 182 Skylane |
| 467 | Swearingen Metro III |
| 085 | Stinson SR-9 |
| 675 | Embraer-145 |
| 050 | Howard DGA-15P |
| 175 | Lockheed L-12A/L-10/10A |
| 406 | Beech 200 Super Kingair |
| 040 | De Havilland DHC2 Beaver |
| 691 | Airbus Industrie A300-600/R/CF/RCF |

- Joining the tables aircraft and flight ordering by trips

```
1  SELECT
2     aircraft,
3     SUM(departures) AS trips
4  FROM flights
5  JOIN aircraft using (aircraft_code)
6  GROUP BY aircraft
7  ORDER BY trips DESC
8  LIMIT 10;
```

Run Query    Explain Query                                            JSON

| aircraft | trips |
|---|---|
| Boeing 737-300 | 19632153 |
| McDonnell Douglas DC9 Super 80/MD81/82/83/88 | 18608868 |
| Canadair RJ-200ER /RJ-440 | 12062704 |
| Boeing 757-200 | 10768257 |
| Boeing 727-200/231A | 9188041 |
| Embraer-145 | 9184729 |
| Boeing 737-100/200 | 8567467 |
| Boeing 737-700/700LR | 7550737 |
| McDonnell Douglas DC-9-30 | 7105295 |
| Airbus Industrie A320-100/200 | 6664803 |

- Analyzing the performance of the queries by fetching the query execution plan.



- Analyzing the data on the table flight

```
1  ANALYZE COMPRESSION flights;
```

Run Query   Explain Query

| Table | Column | Encoding | Est_reduction_pct |
|-------|--------|----------|-------------------|
| flights | year | zstd | 90.86 |
| flights | month | zstd | 88.54 |
| flights | day | zstd | 69.61 |
| flights | carrier | zstd | 49.04 |
| flights | origin | zstd | 46.21 |
| flights | dest | zstd | 40.47 |
| flights | aircraft_code | zstd | 43.32 |
| flights | miles | zstd | 91.08 |
| flights | departures | zstd | 33.30 |
| flights | minutes | zstd | 65.37 |
| flights | seats | zstd | 68.08 |
| flights | passengers | az64 | 0.00 |

- Creating tables from other table

📁 **public**

⌄ Tables (3)

⊞ aircraft

⊞ airports

⊞ flights

```
1  CREATE TABLE airports (
2    airport_code CHAR(3) SORTKEY,
3    airport      varchar(100)
4  );
```

- Copying the data from bucket to the airports table

```
1  COPY airports
2  FROM 's3://us-west-2-aws-training/awsu-spl/spl-17/4.2.5.prod/data/lookup_airports.csv'
3  IAM_ROLE 'arn:aws:iam::723548672068:role/Redshift-Role'
4  IGNOREHEADER 1
5  DELIMITER ','
6  REMOVEQUOTES
7  TRUNCATECOLUMNS
8  REGION 'us-west-2';
```

Run Query   Explain Query

No records found

- Creating a new table vegas_flights to load the data

```
1   CREATE TABLE vegas_flights
2     DISTKEY (origin)
3     SORTKEY (origin)
4   AS
5   SELECT
6     flights.*,
7     airport
8   FROM flights
9   JOIN airports ON origin = airport_code
10  WHERE dest = 'LAS';
```

- Selecting the data loaded from the airport table

```
1  SELECT
2    airport,
3    to_char(SUM(passengers), '999,999,999') as passengers
4  FROM vegas_flights
5  GROUP BY airport
6  ORDER BY SUM(passengers) desc
7  LIMIT 10;
```

Run Query    Explain Query                                                JSON

| airport | passengers |
| --- | --- |
| Los Angeles, CA: Los Angeles International | 29,403,292 |
| Phoenix, AZ: Phoenix Sky Harbor International | 24,160,227 |
| Dallas/Fort Worth, TX: Dallas/Fort Worth International | 15,377,974 |
| Denver, CO: Denver International | 14,937,489 |
| Chicago, IL: Chicago O'Hare International | 14,494,577 |
| San Francisco, CA: San Francisco International | 14,241,188 |
| San Diego, CA: San Diego International | 11,744,708 |
| Salt Lake City, UT: Salt Lake City International | 10,985,774 |
| Atlanta, GA: Hartsfield-Jackson Atlanta International | 10,279,586 |
| Reno, NV: Reno/Tahoe International | 10,166,655 |

- Verifying the disk space and Data distribution using the below query

```
1  SELECT
2    owner AS node,
3    diskno,
4    used,
5    capacity,
6    used/capacity::numeric * 100 as percent_used
7  FROM stv_partitions
8  WHERE host = node
9  ORDER BY 1, 2;
```

Run Query    Explain Query

| node | diskno | used | capacity | percent_used |
| --- | --- | --- | --- | --- |
| 0 | 0 | 1039 | 190633 | 0.5450263070926859400 |
| 1 | 0 | 808 | 190633 | 0.4238510646110589400 |

```
1  SELECT
2    name,
3    count(*)
4  FROM stv_blocklist
5  JOIN (SELECT DISTINCT name, id as tbl from stv_tbl_perm) USING (tbl)
6  GROUP BY name;
```

| Run Query | Explain Query | | JSON | CSV |

| name | | count |
|---|---|---|
| aircraft | ... | 20 |
| airports | ... | 20 |
| flights | ... | 1605 |
| redshift_auto_health_check_1052096 | ... | 16 |
| vegas_flights | ... | 149 |

- Exploring the query performance using redshift console

| Query runtime | Total queries | Workload execution breakdown | Workload concurrency |

| Query ID | Duration | User | SQL | Cluster | Status |
|---|---|---|---|---|---|
| 370,371,372... | 10 s | master | ANALYZE COMPRESSION flights; | ⊘ lab3 | ⊘ Completed |

**Query runtime**



**Cluster status**

● Available    ● Unavailable

## Queries and loads (17)

[Terminate query]

| | Start time ▲ | Query ▼ | Status ▼ | Duration ▼ | SQL ▼ | Copy SQL | PID ▼ | Transaction ID ▼ |
|---|---|---|---|---|---|---|---|---|
| ☐ | Mar 17th, 2020 03:52:02 PM<br>36 minutes ago | 12,37,42 | ⊘ Completed | 2 m | COPY flights FROM 's3://us-west-2 ... | ⧉ Copy | 8710 | 1122 |
| ☐ | Mar 17th, 2020 03:54:01 PM<br>34 minutes ago | 46 | ⊘ Completed | 226 ms | SELECT COUNT(*) FROM flights; | ⧉ Copy | 8710 | 1199 |
| ☐ | Mar 17th, 2020 03:54:57 PM<br>33 minutes ago | 54 | ⊘ Completed | 5 s | SELECT * FROM flights ORDER BY r ... | ⧉ Copy | 8710 | 1215 |
| ☐ | Mar 17th, 2020 03:57:19 PM<br>31 minutes ago | 108 | ⊘ Completed | 1 s | SELECT carrier, SUM (departures) ... | ⧉ Copy | 8710 | 1342 |
| ☐ | Mar 17th, 2020 04:00:23 PM<br>28 minutes ago | 134 | ⊘ Completed | 1 s | SELECT carrier, SUM (passengers) ... | ⧉ Copy | 8710 | 1400 |
| ☐ | Mar 17th, 2020 04:02:54 PM<br>25 minutes ago | 166 | ⊘ Completed | 1 s | SELECT carrier, SUM (miles) FROM ... | ⧉ Copy | 8710 | 1470 |
| ☐ | Mar 17th, 2020 04:06:00 PM<br>22 minutes ago | 198 | ⊘ Completed | 2 s | SELECT carrier, SUM (passengers * ... | ⧉ Copy | 8710 | 1546 |
| ☐ | Mar 17th, 2020 04:08:44 PM<br>20 minutes ago | 226 | ⊘ Completed | 1 s | SELECT carrier, SUM (freight_pou ... | ⧉ Copy | 8710 | 1604 |

- Verifying the query execution details

| ○ | Mar 17th, 2020 03:53:40 PM<br>36 minutes ago | 42 | ⊘ Completed | 2 s | Main | Rewritten query |
|---|---|---|---|---|---|---|

### Query details

| Query ID | Cluster | User | Type | Status |
|---|---|---|---|---|
| 12 | ⊘ lab3 | master | Rewritten query | ⊘ Completed |

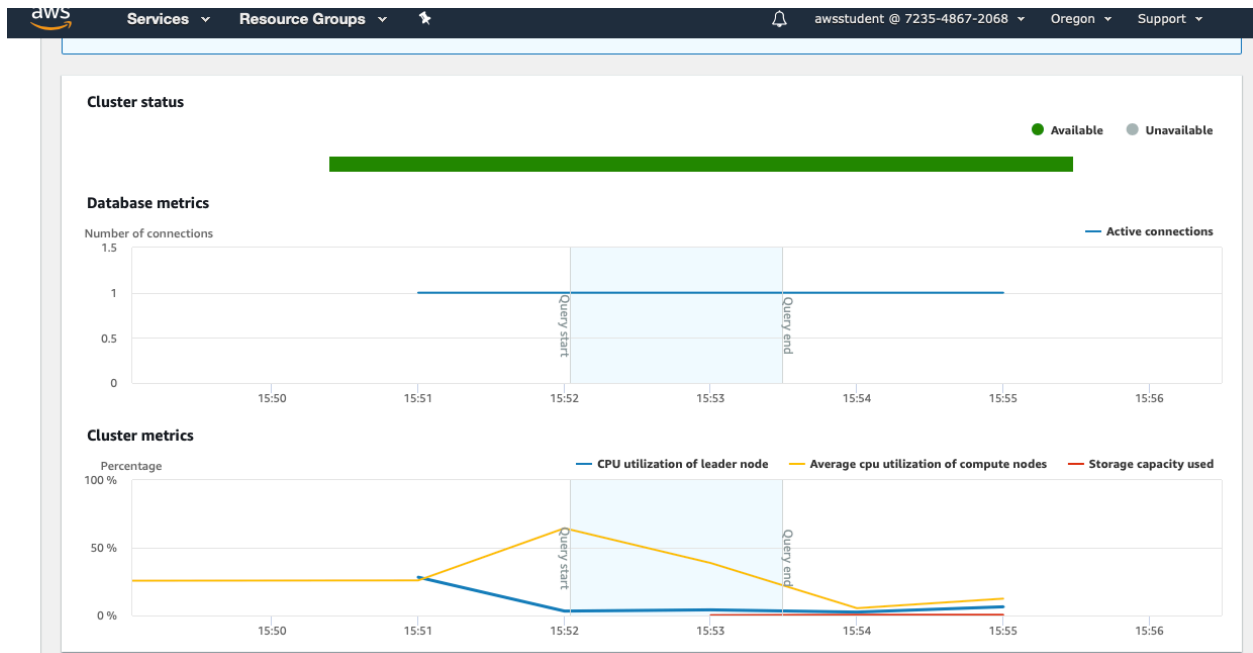| From March 17, 2020 at 03:52:02 PM To March 17, 2020 at 03:53:29 PM | Total runtime | 1m |
|---|---|---|

**Query details** | Query plan

**SQL** [Copy]

```
COPY flights
FROM
    's3://us-west-2-aws-training/awsu-spl/spl-17/4.2.5.prod/data/flights-usa' IAM_ROLE '' GZIP
DELIMITER ',' REMOVEQUOTES REGION 'us-west-2';
```

**Execution details**

| Total runtime | 1m |
|---|---|
| **Statistics** | |
| Rows returned | 484,128,783 |
| Total data scanned | 21.39 GB |

- Verifying the database metrics using cluster console

- Verifying the database metrics using cluster console



- Deleting the cluster after use.