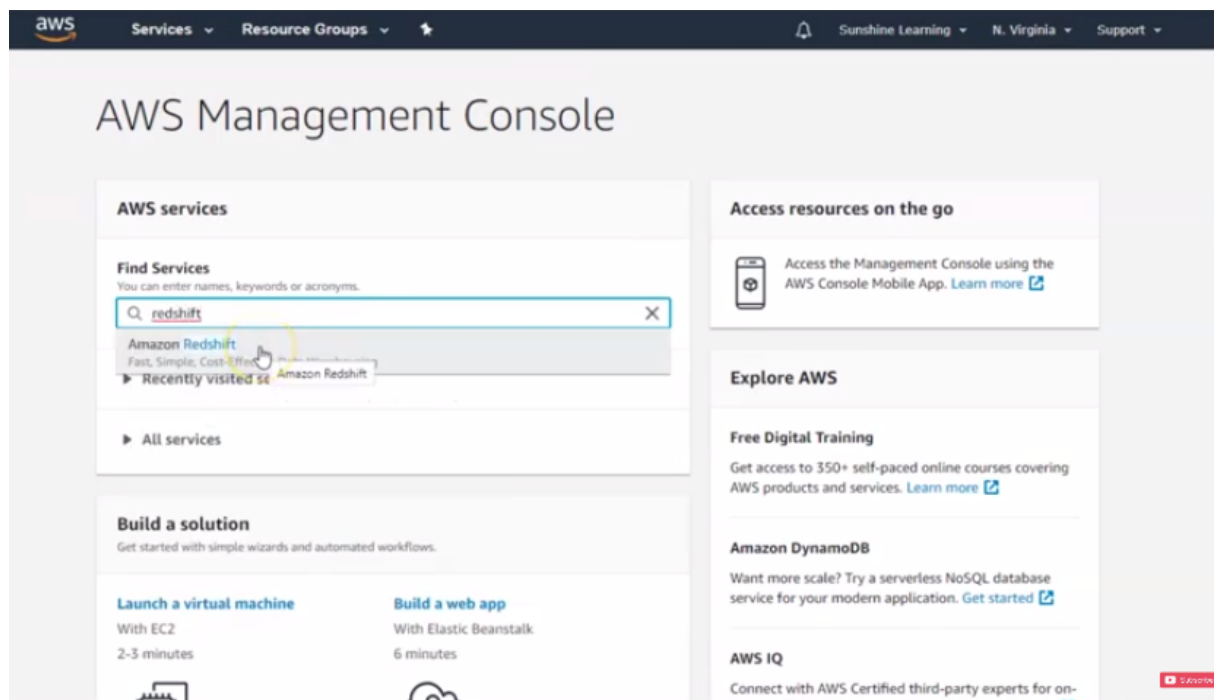# CS 79B Final Project - Part 1

Identify an <u>AWS service</u> from the <u>Database, Analytics or Migration category group</u> not discussed in class from the AWS Console. Make sure that the underlined items are covered in your answer. Create a word or text document and upload it.

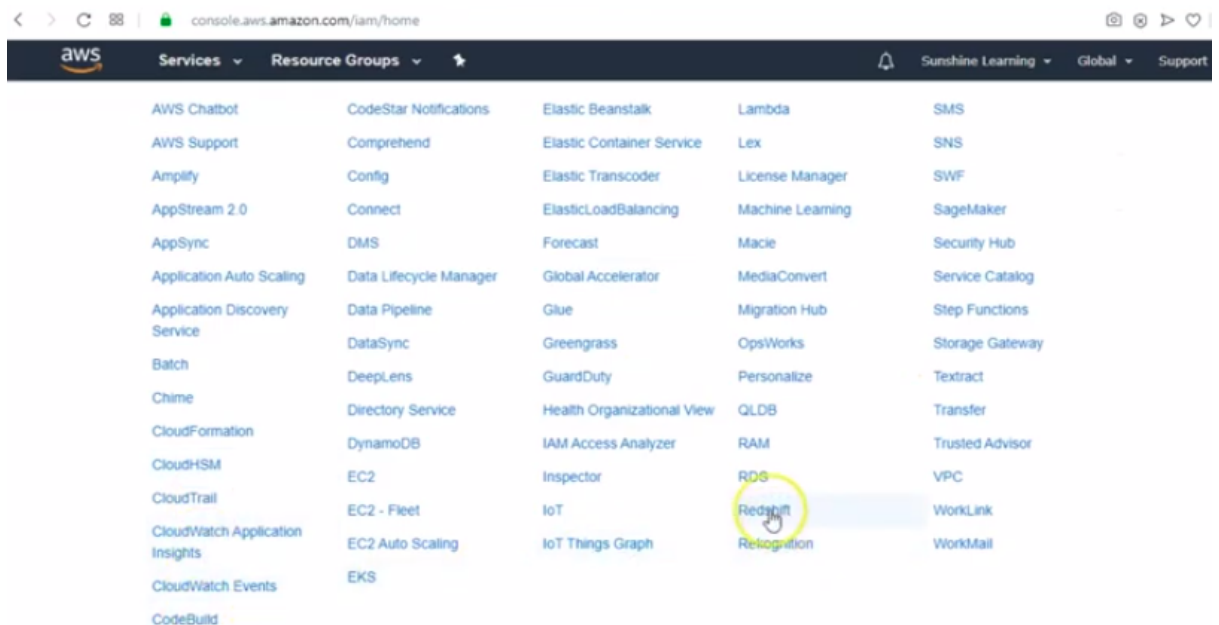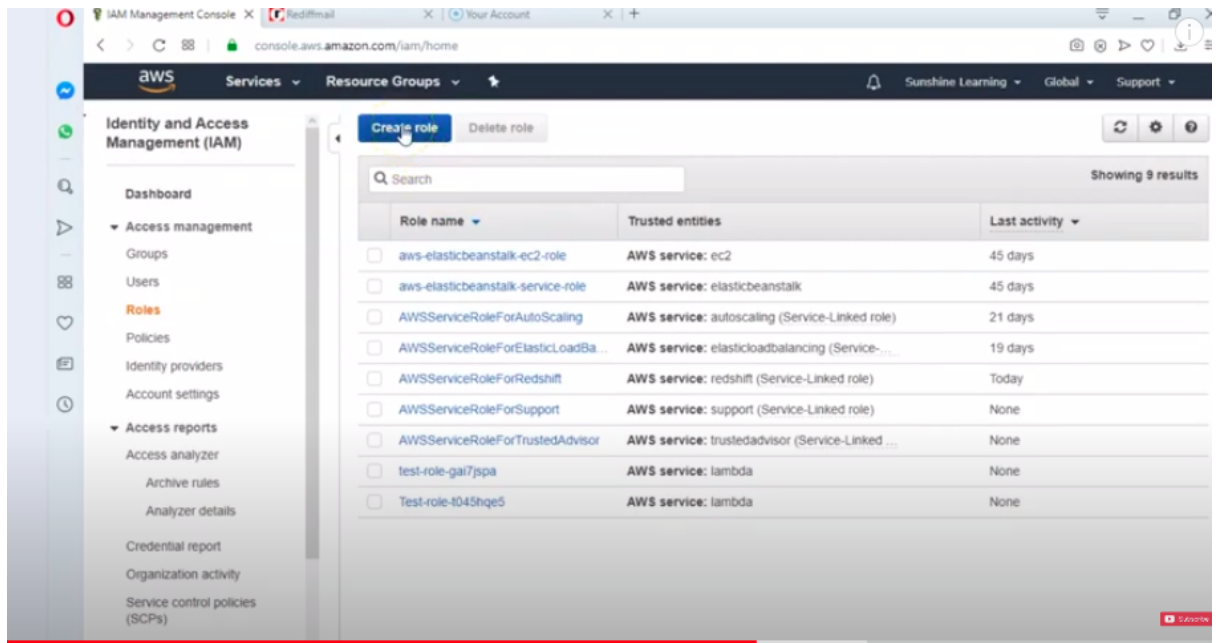**The AWS Service selected for the Final Project is AWS Redshift**

(1) For this service, <u>describe the typical way to interact with the service</u> (that is, a list of actions or steps users' takes to implement this service. Similar to instructions provided to create a MySQL DB, Image Rekognition etc.)

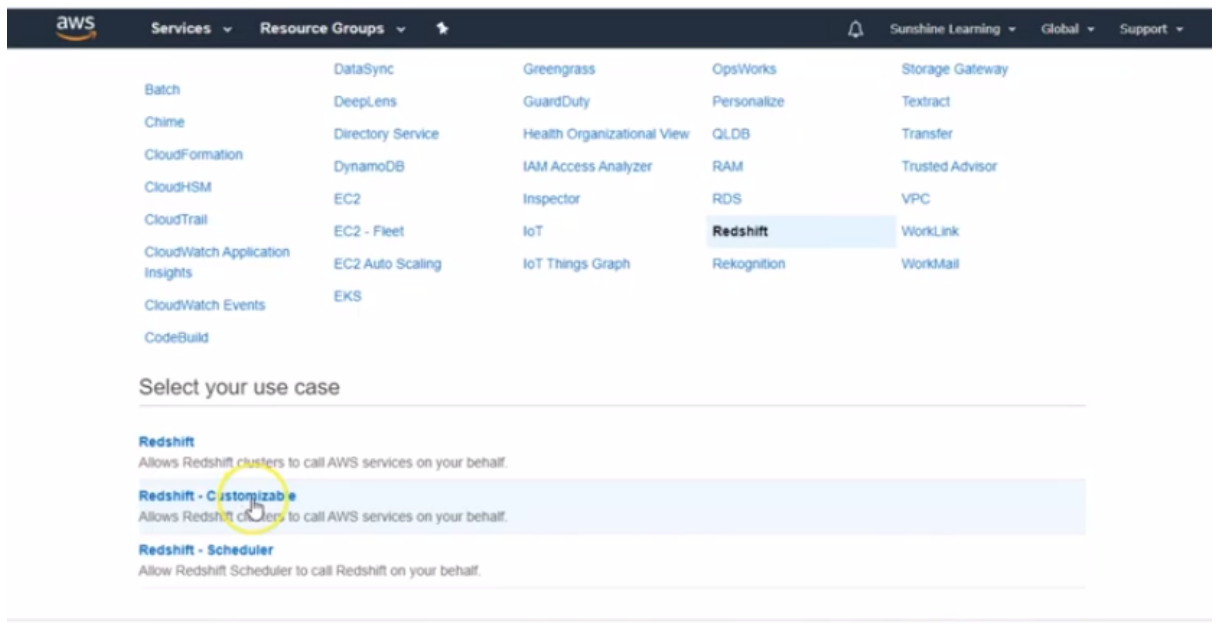**Steps to interact with Amazon Redshift:**

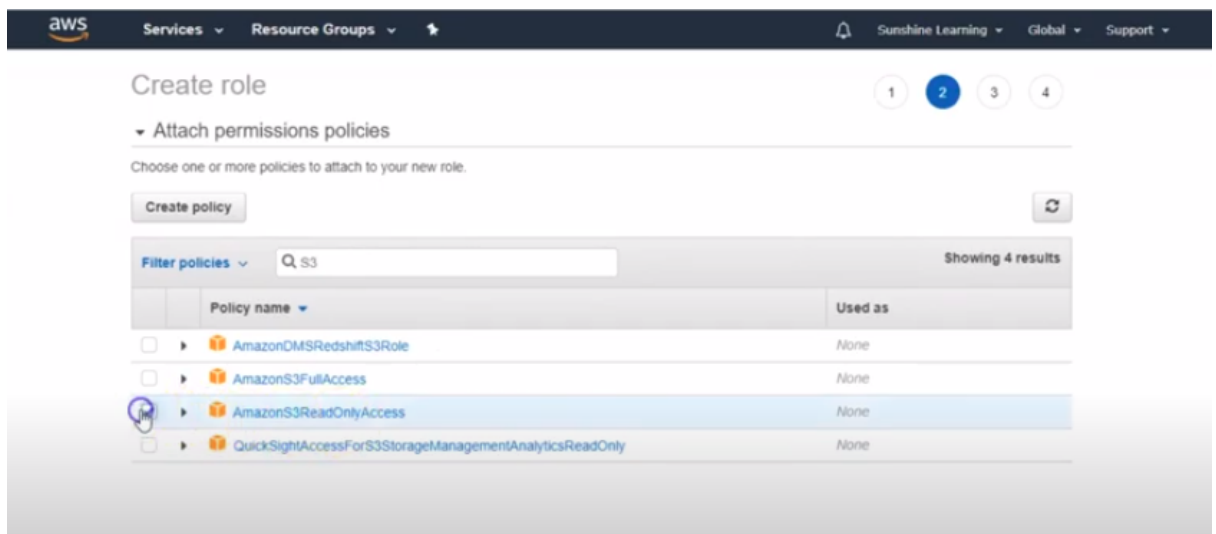1. *Search for Amazon Redshift in Amazon Management Console*



2. *Before creating Cluster, we need to **create an IAM role**. Select **Redshift** under IAM Role because Redshift will be calling S3 service. We will putting up data in S3 and then it will be uploaded to Redshift*

3. Click on **Redshift Customizable** in the Use Case

4. Under Policy Permissions, select **AMAZONS3READONLYACCESS** and assign permissions to this Role



5. Enter **Role** name and click on **Create Role**

6. *Next, go to **Redshift** and create a **Cluster**. Also select the **uncompressed data size**: GB/TB/PB*

7. Give **Cluster name, No: of Nodes, Master Username and Password** for Redshift Cluster

8. *You can change **Default VPC** from Default Settings*



9. **Redshift Cluster** *is being created*

10. To run SQL queries, you can use **Query Editor**

11. To create Tables in Redshift, use **SQL Workbench/J** client to connect to **Redshift DB**



12. Connect to Redshift DB using **JDBC URL** with **PORT 5439.** Select the **Driver, Username and Password**

13. Download sample data from below location and load it to **S3**

14. Write queries to **Create** table, Use **copy** statement to copy file from S3 to Redshift DB, Use **Select** statement to query the results from Sales table

*15. In **copy** statement **credentials**, define the **Role ARM***



*16. Execute all these 3 queries. You should be able to Connect to Database and display the tables in Amazon Redshift*

(2) How does AWS charge for this service?

Amazon Redshift costs less to operate than any other data warehouse. It starts small at $0.25 per hour and scale up to petabytes of data and thousands of concurrent users.

**On-Demand Pricing:**

Example:

In US West (Northern California), for Dense Compute EC2 dc2.large -  $0.33 per Hour

**Redshift Pricing Spectrum:**

$5.00 per terabyte of data scanned

**Concurrency Scaling pricing:**

A 10 DC2.8XL node Redshift cluster in the US-East costs $48 per hour.

**Redshift managed storage pricing:**

In US West (Northern California), for Storage/month - $0.0271 per GB


(3) Is it a part of the free tier?

The Amazon Redshift free trial program is not part of the AWS Free Tier

(4) Provide URLs to documentation for this service.

https://aws.amazon.com/redshift/

https://docs.aws.amazon.com/redshift/latest/dg/c_high_level_system_architecture.html

https://aws.amazon.com/redshift/pricing/

https://docs.aws.amazon.com/redshift/latest/dg/welcome.html


(5) URLs to AWS videos about this service.

https://www.youtube.com/watch?v=7bfOllAyxlg


(6) Watch the videos you located.

https://www.youtube.com/watch?v=7bfOllAyxlg


(7) For one of the videos you listed, answer the following questions:

- **Summarize a few key points made in the video.**

**Definition of Amazon Redshift:** It's a Cloud-based Data warehouse service that is used for collecting and storing data. It's also used to get and analyze the data using BI (Business Intelligence) tools and simplifies the process of handling large sets of data.

**Data warehouse Definition:** Data warehouse is a repository where the data is stored

1. Before Redshift, people used to fetch the data from the Data warehouse.
2. Fetching data from Data warehouse was a complicated task because the developer and Data warehouse might be located in different geographic locations and there might be network connectivity issues, internet connectivity challenges, security challenges, and a lot of maintenance is required to manage Data warehouse.



3. **Cons of Traditional Data warehouse service:**
   a. It's time consuming process to download or get data from Data warehouse
   b. Maintenance cost was high
   c. There is possibility of losing information in between downloading of data
   d. Data rigidity was an issue
4. All these problems can be solved with Amazon Redshift
5. **Pros:**
   a. **Costs** less when compared to other Cloud data warehouse products in the market. You can have a large Data warehouse or combine databases in a Data warehouse at a very low cost.
   b. It's the fastest Data warehouse (in performance) in market
   c. It has more than 15k customers
   d. Whatever you use, you pay for it ONLY

e. **Scalability:** If you want to increase the nodes in your database, then you can increase it on the fly. You need not switch off to scale.
f. **Availability:** It's highly available in multiple Availability Zones.
g. You can have multiple clusters in Redshift, you can define own VPC and Security Group for each Cluster and increase the security of Redshift
h. **Flexibility:** You can remove cluster, create new cluster, take a snapshot of the deleting cluster, move the cluster to different region
i. You can have a simple migration from a traditional database to a Cloud Redshift data migration by in-built tools in Cloud being connected to a traditional database.

● **Identify two interesting quotes that were made.**

   **Architecture of Amazon Redshift:**
   ○ We have Compute Node which does data processing
   ○ Leader Node which gives instructions to Compute Node and it also manages Client applications that require data from Redshift.
   ○ The client applications connect with JDBC (Java Database Connectivity) and ODBC (Open Database Connectivity).
   ○ Amazon Redshift can monitor connections from Client applications using JDBC.
   ○ ODBC allows Client applications to have live or direct data interaction with Amazon Redshift. The Leader Node can get information from Compute nodes.
   ○ Compute nodes (Nodes) processes the data, they are a set of computing resources and when combined together are called Clusters.
   ○ A Cluster is a set of Computing nodes called Nodes.
   ○ Nodes which are combined together are called Data warehouse Cluster.
   ○ We can have 1...100 Compute nodes which is a scalable solution.
   ○ Each Cluster has a Database in the form of Node.
   ○ The Leader Node manages the interaction between Client application and Compute Node i.e acts as bridge between both. It also analyzes and develops designs in order to carry out database operations.
   ○ At high-level, Leader Node sends out instructions to be performed or executed by Compute Node and sends the response to Client applications.
      ■ The Leander Node runs the programs and assigns the code to individual compute nodes.
      ■ The Compute Node executes the program and shares the results to Leader Node for final aggregation and then it's delivered to Client application
   ○ Each Computer Node is categorized into slices and each slice is allocated with specific memory space, where it processes it's workload.
   ○ These Node slices work in parallel and hence the reason why Redshift is the fastest Data warehouse when compared to traditional and currently existing Data warehouses

● **What new facts did you learn from watching this video.**

1. We have additional concepts called Column storage and Compression.
2. **Column storage:** Data is stored in columns which helps in optimizing query performance and quicker output. Below is the example. It makes the data more structured and easy to extract.
3. Compression: To save Column storage we can use Compression as an attribute. It's a Column-level operation which decreases storage requirements and improves query performance.

● **What was the best part of the video?  Why?**

   **Use cases:** DNA, a Telecommunications company is facing issues with handling Website data and Amazon S3 data which leads to slow process of the application. They overcame this issue by using Amazon Redshift and they noticed a 52% increase in application performance.

● **What questions remain in your mind after watching the video?  Why?**

 **None**