# Project Deliverable 2

**Team 10:**
**Sanjana**
**Tarun**
**Swathi**
**Archana**
**Shashank**

**Set Up AWS Environment - Account Setup: Team Access to AWS Learner Lab and Github, create IAM Roles**
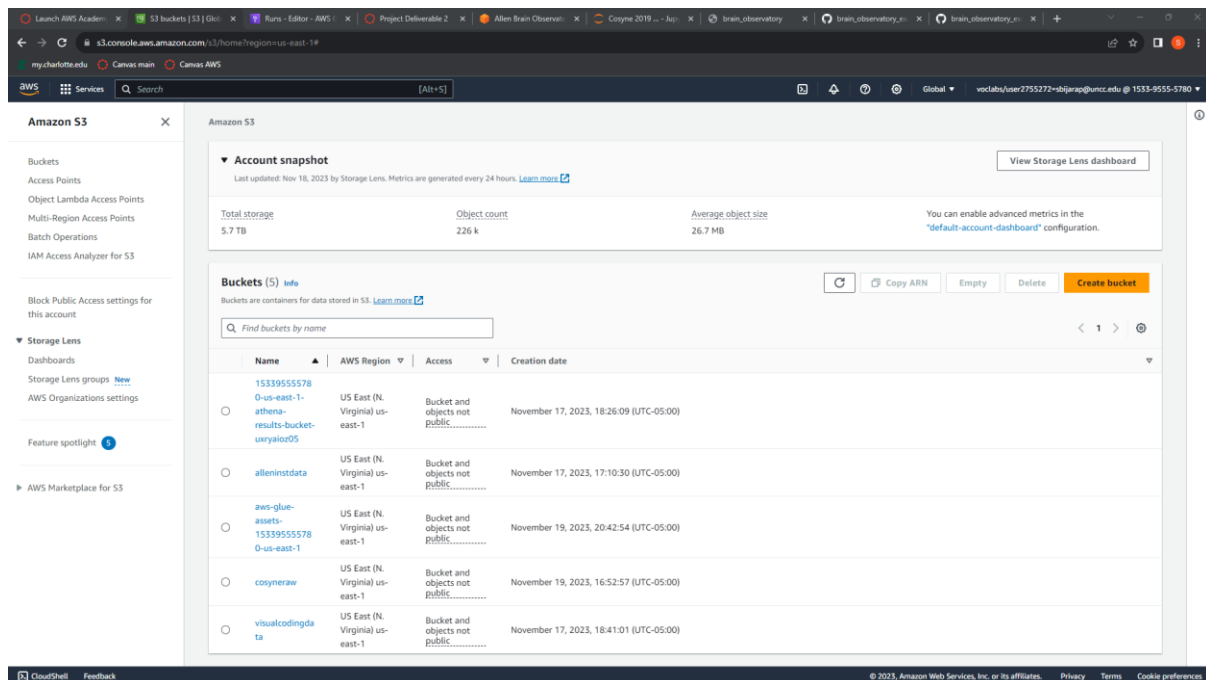
https://github.com/RoshnaSanjanaKommareddy/Cloud_Project_10#cloud_project_10.
The whole team had access to AWS Learner Lab and IAM roles were created.

**S3 Data Storage:**
**Create an S3 bucket to store the dataset(s)**
**Upload data to the S3 Bucket**



**Data Exploration for Insight and Pre-processing**
**Use Amazon Athena to query the transformed data.**
**Use SQL queries for meaningful insights from the dataset for data exploration**
**Create visualizations using Amazon QuickSight.**

**Query 3 | Query 4**

```sql
1  CREATE TABLE visualcodingdatabase.cosney_processed
2  WITH (
3      external_location = 's3://cosyneraw/cosney_processed'
4  ) AS
5  SELECT col1 AS cell_id,
6      col2 AS cre,
7      col3 AS area,
8      col4 AS expt_id,
9      col5 AS signal_correlation_dg,
10     col6 AS signal_correlation_nm3,
11     col11 AS number_cell,
12     col18 AS imaging_depth,
13     col9 AS age,
14     col10 AS sex
15     FROM visualcodingdatabase.cosyneraw;
```

SQL  Ln 3, Col 55

Run again  Explain  Cancel  Clear  Create

Reuse query results
up to 60 minutes ago

Query results  Query stats

⊘ Completed    Time in queue: 58 ms    Run time: 1.578 sec    Data scanned: 11.21 KB

Query successful.

---

**Query 3 | Query 4 | Query 5**

```sql
1  select * from cosney_processed
```

SQL  Ln 1, Col 31

Run again  Explain  Cancel  Clear  Create

Reuse query results
up to 60 minutes ago

Query results  Query stats

⊘ Completed    Time in queue: 157 ms    Run time: 1.033 sec    Data scanned: 5.17 KB
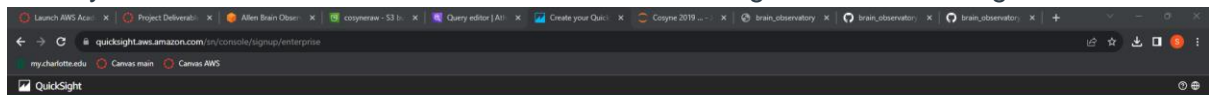
Results (117)    Copy    Download results

| # | cell_id | cre | area | expt_id | signal_correlation_dg | signal_correlation_nm3 | number_cell | imaging_depth | age | sex |
|---|---------|-----|------|---------|----------------------|------------------------|-------------|---------------|-----|-----|
| 1 | cell_id | Cre | area | expt_id | signal_correlation_dg | signal_correlation_nm3 | numbercells | imaging_depth | age | sex |
| 2 | 598998745 | Sst | VISl | 597028936 | 0.35121026913019643 | 0.0745393538348721 | 12.0 | 375.0 | 105.0 | female |
| 3 | 662056046 | Sst | VISp | 612044633 | 0.3832542149008148 | 0.1397119205998549 | 8.0 | 265.0 | 98.0 | male |
| 4 | 662098491 | Sst | VISpm | 639117194 | 0.11593636022081612 | 0.11518353728960953 | 7.0 | 375.0 | 98.0 | male |

Amazon Quicksight asked for a subscription account and even after accepting to pay it says that my role as the learner lab user restricts me from using Amazon QuickSIght



**AWS Glue ETL Job:**
**Create an ETL job using AWS Glue to transform the dataset(s) in S3**
**Perform basic data transformations (e.g., filtering, aggregation, type conversions)**
**AWS Pipeline/Solution Chart**