# Report on Superstore Dataset Analysis

## OBJECTIVE:

Analysis of a Superstore

## ABSTRACT OF THE PROJECT:

The project is to develop Analysis of a Superstore.
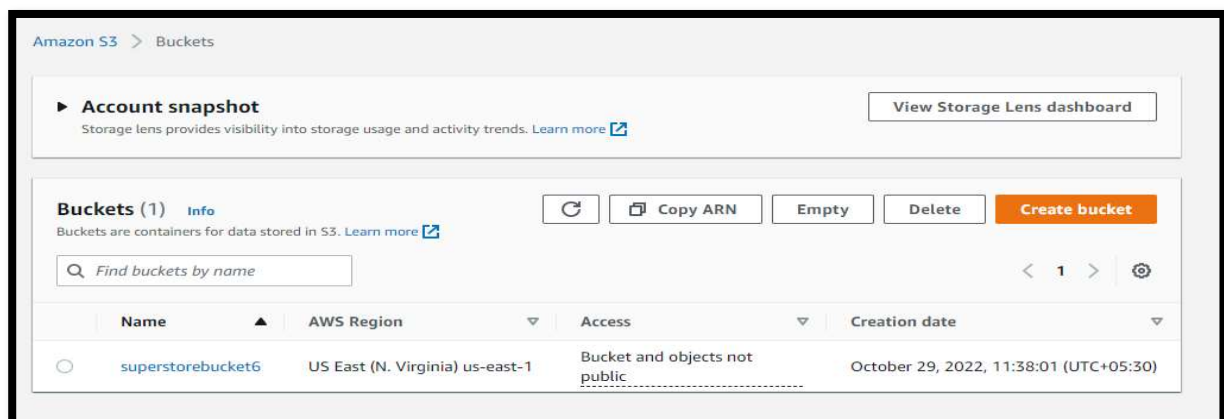
## TECHNOLOGY USED:

- Snowflake
- Cloud (AWS/Azure)
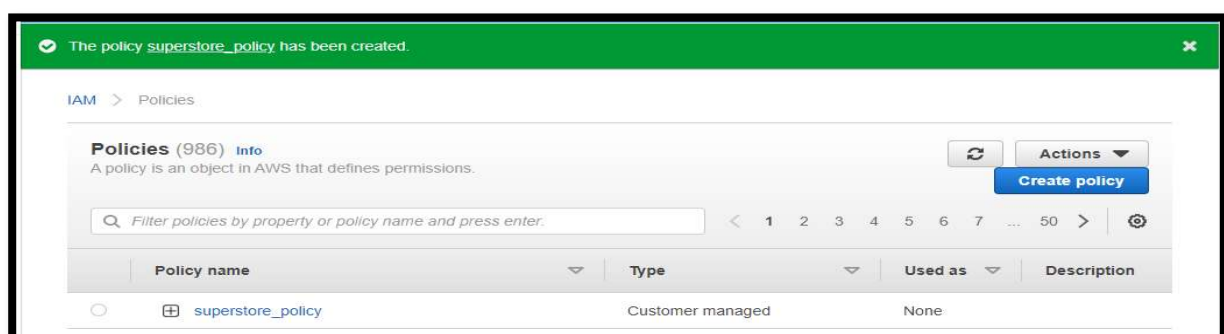
## IMPLEMENTATION:

### 1. Create an External Stage to load the data continuously
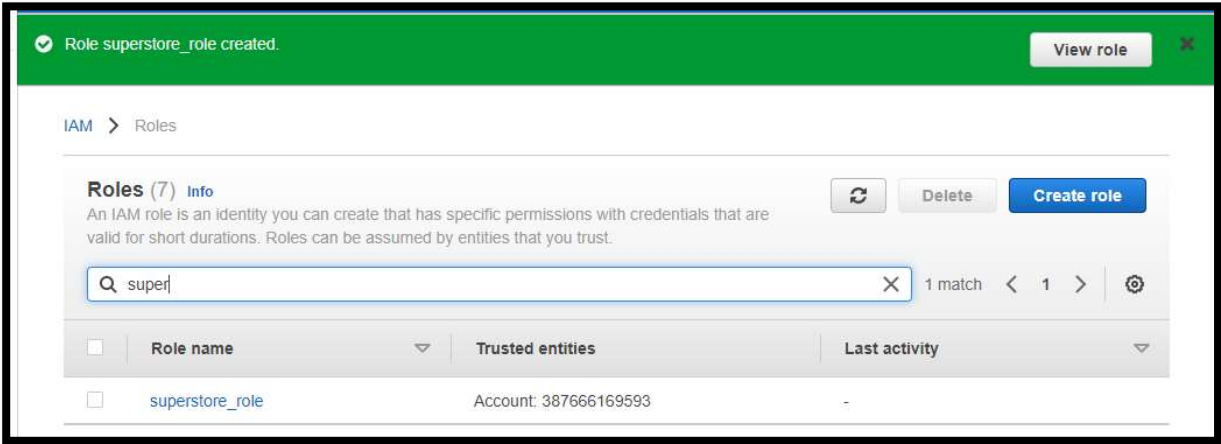
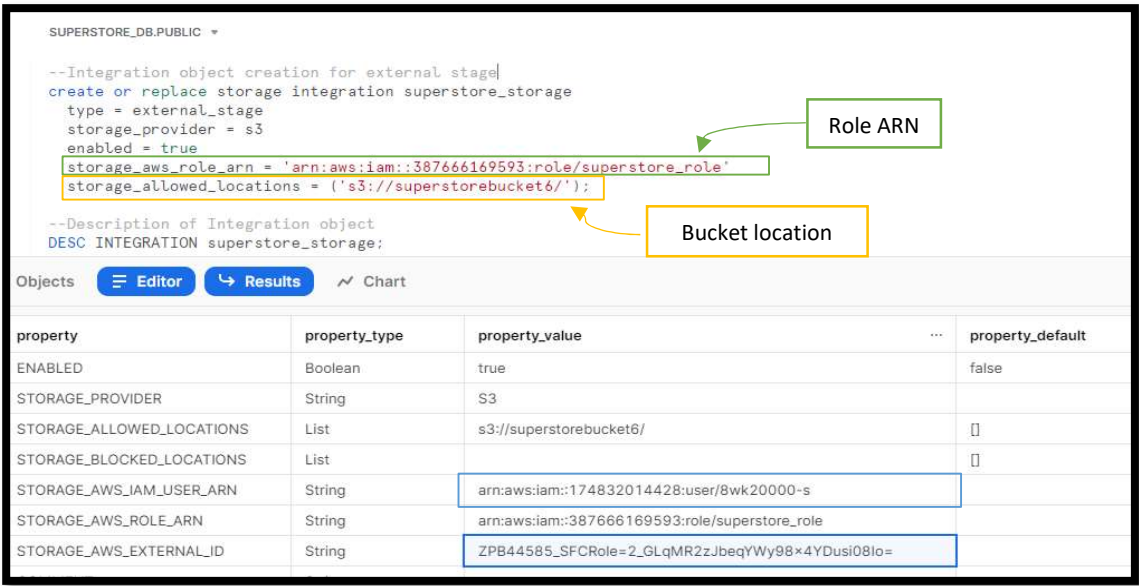In AWS,

Step 1: Create Bucket
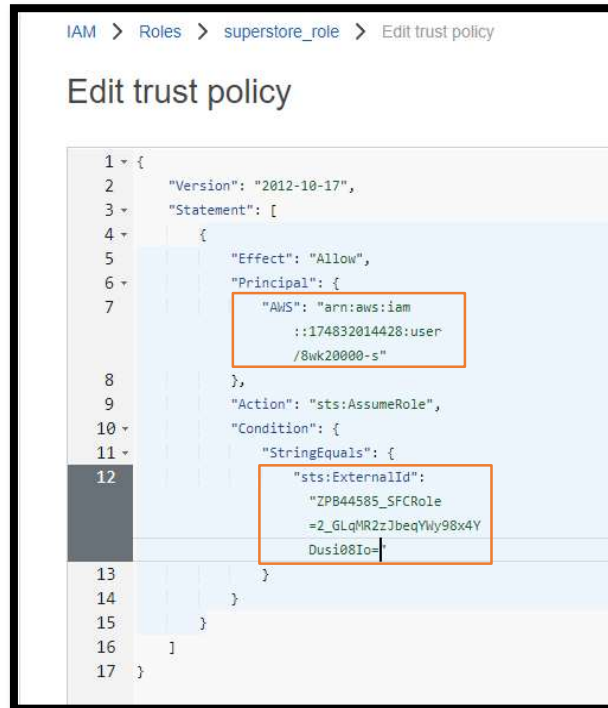


Step 2: Create Policy

## Step 3: Create Role



In Snowflake,

## Step 4: Integration Creation (Integrating the bucket with storage)

## Step 5: Editing trust policy



## Step 6: Create External stage and listing the files

**Load the data from the External Stage to the respective Table**

Step 1: Created Table superstore_table

Step 2: Loading data into table superstore_table from external stage

## 2. Create a scheduler to schedule the job at 12:00 AM IST hours every Thursday to perform previous step

Step 1: Created a task to schedule a job at 12:00 AM IST hours every Thursday

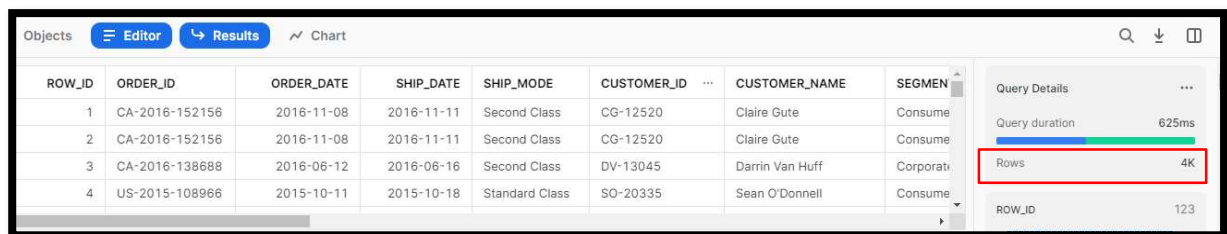Step 2: The task will load the data from external stage to source table

Before task was executed:



After task was executed:

**3. Create an alternative of the above step for auto ingestion of new record sets uploaded on the blob storage into respective snowflake table.**

Step 1: Snowpipe created with auto-ingestion

```
create or replace pipe superstore_pipe auto_ingest=true as
    copy into superstore_table
    from @superstore_stage;
```

Step 2: Description of snowpipe



Step 3: Event Notification created in AWS S3



Before Upload Source table has 4k rows :



Step 4: After Uploading source file into S3, Fetching all data from source table

## 4. Implement SCD 2 in for the above data flow

Step 1: Created a stream to capture DML changes made to the source table

Step 2: Created a task to merge data into target table

```
--Merged data into target table with versioning history which will capture through task
CREATE OR REPLACE TASK super_target_merge
  WAREHOUSE = superstore_wh
  SCHEDULE = '1 minute'
WHEN
  SYSTEM$STREAM_HAS_DATA('superstore_stream')
AS
merge into super_target t
using superstore_stream s
on t.Row_ID=s.Row_ID and (metadata$action='DELETE')
when matched and metadata$isupdate='FALSE' then update set rec_version=9999, stream_type='DELETE'
when matched and metadata$isupdate='TRUE' then update set rec_version=rec_version-1, stream_type='UPDATE'
when not matched then insert
  (Row_ID,Order_ID,Order_Date,Ship_Date,Ship_Mode,Customer_ID,Customer_Name,Segment,Country,City,State,Postal_Code,Region,Product_ID,Category,
                    Sub_Category,Product_Name,Sales,Quantity,Discount,Profit,stream_type,rec_version,REC_DATE)
  values(s.Row_ID,s.Order_ID,s.Order_Date,s.Ship_Date,s.Ship_Mode,s.Customer_ID,s.Customer_Name,s.Segment,s.Country,s.City,s.State,s.Postal_Code,

  s.Region,s.Product_ID,s.Category,s.Sub_Category,s.Product_Name,s.Sales,s.Quantity,s.Discount,s.Profit, metadata$action,0,CURRENT_TIMESTAMP());
```

Step 3: Displayed the target table that will return all the records with versioning history

Update:

```
update superstore_table set order_id='CA-2022-DDDDDD' where row_id=777;
select * from super_target where row_id=777;
```

| | SALES | QUANTITY | DISCOUNT | PROFIT | STREAM_TYPE ... | REC_VERSION | REC_DATE |
|---|---|---|---|---|---|---|---|
| Pencil, #2 | 32.76 | 7 | 0.2 | 3.6855 | INSERT | 0 | 2022-10-29 23:26:26.874 -0700 |
| Pencil, #2 | 32.76 | 7 | 0.2 | 3.6855 | UPDATE | -1 | 2022-10-29 23:24:34.222 -0700 |

Delete:

```
delete from superstore_table where row_id=777;
select * from super_target where row_id=777;
```

| | SALES | QUANTITY | DISCOUNT | PROFIT | STREAM_TYPE ... | REC_VERSION | REC_DATE |
|---|---|---|---|---|---|---|---|
| Pencil, #2 | 32.76 | 7 | 0.2 | 3.6855 | DELETE | 9,999 | 2022-10-29 23:26:26.874 -0700 |

**5. Implement Row Level and Column Level Security on the above dataset.**

1. Column level Security

- Created a new role (superstore_role)
- Created masking policy
- Masked sensitive columns from target table when viewed from any other   role.

Viewing target table through superstore_role:
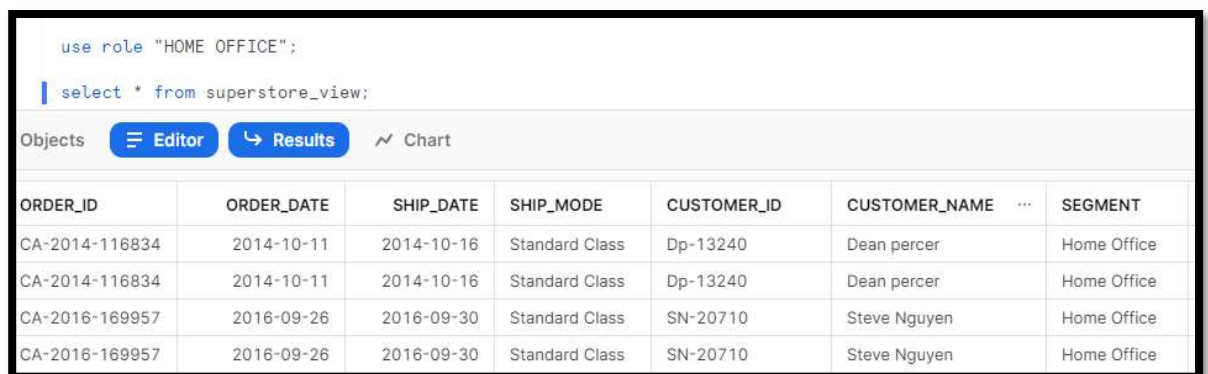


Viewing target table through other role:



2. Row level Security:

- Created roles based on a relevant column
- Created users with their default roles
- Granted secure view of respective entries/rows to the roles.

Viewing superstore data for Consumer role:



Viewing superstore data for Home Office user:



Negative testing :