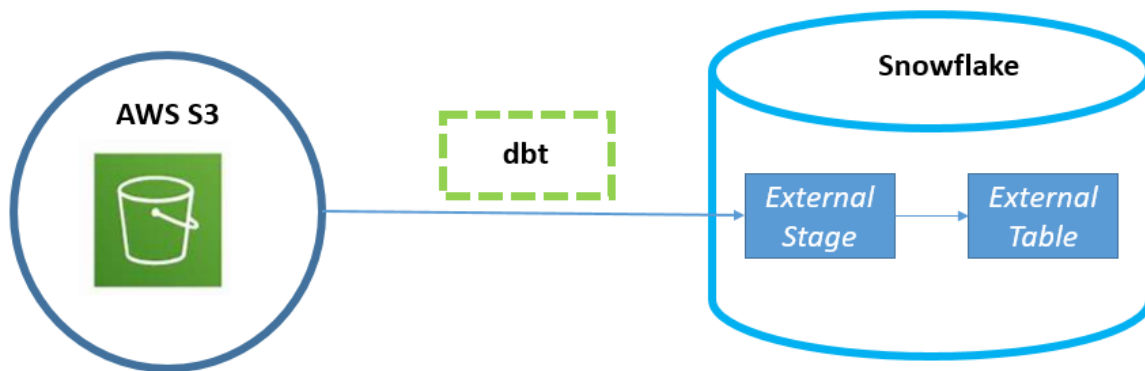


Snowflake with dbt(Data Build Tool): ETL Transformations (Installed via CLI)

*I designed and implemented an end-to-end ETL pipeline leveraging **AWS S3**, **Snowflake**, and **dbt (through CLI)** to build a scalable, structured data workflow. The project included orchestrating data ingestion from S3 into Snowflake, performing advanced common table expressions(**CTE's**) of **SQL** transformations using dbt, and managing clean, version-controlled development through GitHub. This hands-on build provided valuable practical insights into cloud-based data engineering and modern ELT practices, resulting in a robust pipeline aligned with industry standards.*



"How I built an end-to-end Snowflake + dbt pipeline from scratch 🚀"

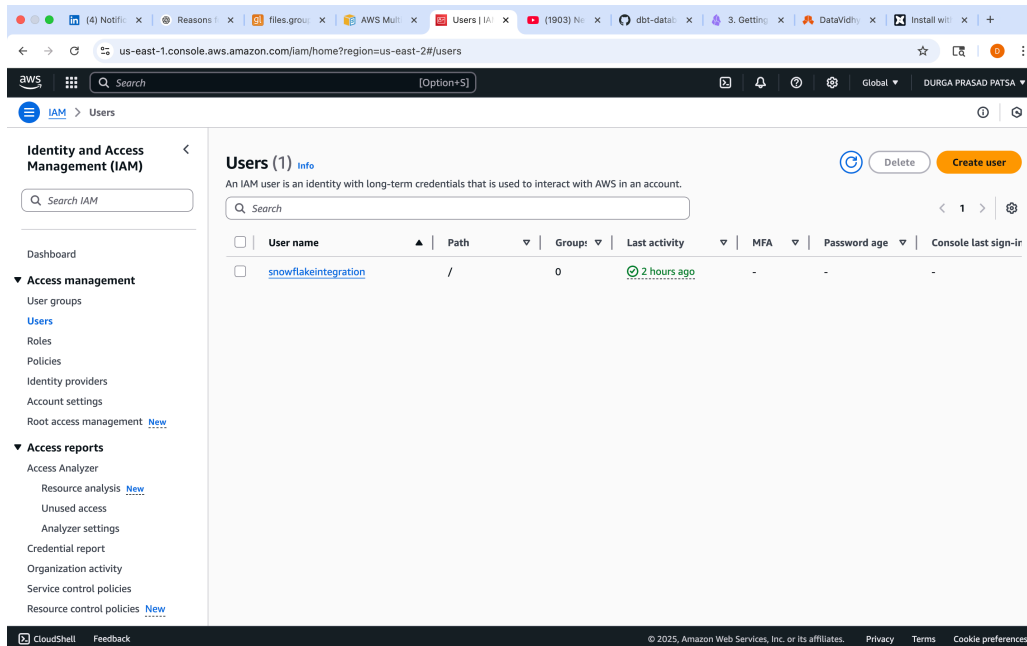
- 1 Downloaded MovieLens CSV files locally.
- 2 Uploaded the CSVs to an S3 bucket in my AWS workspace.
- 3 Created an IAM role with access keys for Snowflake to securely pull data.
- 4 Set up Snowflake: created users, roles, warehouse, database, schema.
- 5 Pulled CSV data from S3 into Snowflake raw tables using CREATE STAGE and COPY INTO commands.
- 6 Installed dbt CLI (instead of dbt Cloud) and initialized a structured dbt project in VS Code.

- 7 Created **three folders**: staging, dim, fct for organized Medallion architecture.
 - 8 Built src_* models in staging as **views** from raw Snowflake tables.
 - 9 Transformed src_* views to dim_* tables with cleaning and structuring logic.
 - 10 Further transformed dim_* tables to fct_* fact tables to create a clean, analytics-ready final table.
- ✓ Validated the pipeline with dbt run, dbt test, and explored the lineage in dbt docs.

S3 Work: Upload CSV files into AWS S3 and setup the IAM role and access.

The screenshot shows the AWS S3 console interface. The left sidebar displays the 'Amazon S3' navigation menu with options like 'General purpose buckets', 'Directory buckets', 'Table buckets', 'Access Grants', 'Access Points for general purpose buckets', 'Access Points for directory buckets', 'Object Lambda Access Points', 'Multi-Region Access Points', 'Batch Operations', 'IAM Access Analyzer for S3', 'Block Public Access settings for this account', 'Storage Lens', 'Dashboards', 'Storage Lens groups', 'AWS Organizations settings', and 'Feature spotlight'. The main content area is titled 'netflixpatsa' and includes tabs for 'Objects', 'Metadata', 'Properties', 'Permissions', 'Metrics', 'Management', and 'Access Points'. The 'Objects' tab is active, showing a list of 6 objects. Above the list are buttons for 'Copy S3 URI', 'Copy URL', 'Download', 'Open', 'Delete', 'Actions', 'Create folder', and 'Upload'. A search bar is present with the text 'Find objects by prefix'. The object list has columns for 'Name', 'Type', 'Last modified', 'Size', and 'Storage class'. The objects listed are: 'genome-scores.csv' (308.6 MB), 'genome-tags.csv' (17.7 KB), 'links.csv' (556.7 KB), 'movies.csv' (1.3 MB), 'ratings.csv' (508.7 MB), and 'tags.csv' (15.8 MB). All objects are of type 'csv' and 'Standard' storage class.

Name	Type	Last modified	Size	Storage class
genome-scores.csv	csv	July 9, 2025, 10:15:01 (UTC-05:00)	308.6 MB	Standard
genome-tags.csv	csv	July 9, 2025, 10:15:01 (UTC-05:00)	17.7 KB	Standard
links.csv	csv	July 9, 2025, 10:15:01 (UTC-05:00)	556.7 KB	Standard
movies.csv	csv	July 9, 2025, 10:15:02 (UTC-05:00)	1.3 MB	Standard
ratings.csv	csv	July 9, 2025, 10:15:01 (UTC-05:00)	508.7 MB	Standard
tags.csv	csv	July 9, 2025, 10:15:04 (UTC-05:00)	15.8 MB	Standard



1 Environment Setup in Snowflake

Create Warehouse, Database, Schema

```
USE ROLE ACCOUNTADMIN;  
CREATE ROLE IF NOT EXISTS TRANSFORM;  
GRANT ROLE TRANSFORM TO ROLE ACCOUNTADMIN;  
  
CREATE WAREHOUSE IF NOT EXISTS COMPUTE_WH;  
GRANT OPERATE ON WAREHOUSE COMPUTE_WH TO ROLE TRANSFORM;  
  
CREATE DATABASE IF NOT EXISTS MOVIELENS;  
CREATE SCHEMA IF NOT EXISTS MOVIELENS.RAW;  
  
GRANT ALL ON WAREHOUSE COMPUTE_WH TO ROLE TRANSFORM;  
GRANT ALL ON DATABASE MOVIELENS TO ROLE TRANSFORM;  
GRANT ALL ON ALL SCHEMAS IN DATABASE MOVIELENS TO ROLE TRANSFORM;
```

Create User for dbt

```
CREATE USER IF NOT EXISTS dbt PASSWORD = '*****'  
DEFAULT_WAREHOUSE = 'COMPUTE_WH'  
DEFAULT_ROLE = TRANSFORM  
DEFAULT_NAMESPACE = 'MOVIELENS.RAW';  
  
GRANT ROLE TRANSFORM TO USER dbt;
```

Create External Stage for S3 Ingestion

```
CREATE STAGE netflixstage
URL = 's3://netflixpatsa'
CREDENTIALS = (AWS_KEY_ID='***' AWS_SECRET_KEY='***');
```

2 Ingesting CSV Data from S3 into Snowflake Staging Tables

Create staging tables and load data:

#Creating raw movies table and copying data from S3 location

```
CREATE OR REPLACE TABLE raw_movies (

    movieId INTEGER,
    title STRING,
    genres STRING

);
COPY INTO raw_movies FROM '@netflixstage/movies.csv' FILE_FORMAT = (TYPE =
'CSV' SKIP_HEADER = 1 FIELD_OPTIONALLY_ENCLOSED_BY = '');
```

Creating raw ratings table and copying data from S3 location

```
CREATE OR REPLACE TABLE raw_ratings (
    userId INTEGER,
    movieId INTEGER,
    rating FLOAT,
    timestamp BIGINT
);
COPY INTO raw_ratings FROM '@netflixstage/ratings.csv' FILE_FORMAT = (TYPE =
'CSV' SKIP_HEADER = 1 FIELD_OPTIONALLY_ENCLOSED_BY = '');
```

Creating raw movies tags and copying data from S3 location

```
CREATE OR REPLACE TABLE raw_tags (
    userId INTEGER,
    movieId INTEGER,
    tag STRING,
    timestamp BIGINT
);
COPY INTO raw_tags FROM '@netflixstage/tags.csv' FILE_FORMAT = (TYPE = 'CSV'
SKIP_HEADER = 1 FIELD_OPTIONALLY_ENCLOSED_BY = '') ON_ERROR = 'CONTINUE';
```

Creating raw genome table and copying data from S3 location

```
CREATE OR REPLACE TABLE raw_genome_tags (
    tagId INTEGER,
    tag STRING
);
COPY INTO raw_genome_tags FROM '@netflixstage/genome-tags.csv' FILE_FORMAT =
(TYPE = 'CSV' SKIP_HEADER = 1 FIELD_OPTIONALLY_ENCLOSED_BY = '');
```

Creating raw genome scores table and copying data from S3 location

```
CREATE OR REPLACE TABLE raw_genome_scores (  
    movieId INTEGER,  
    tagId INTEGER,  
    relevance FLOAT  
);  
COPY INTO raw_genome_scores FROM '@netflixstage/genome-scores.csv'  
FILE_FORMAT = (TYPE = 'CSV' SKIP_HEADER = 1 FIELD_OPTIONALLY_ENCLOSED_BY =  
'');
```

Creating raw links table and copying data from S3 location

```
CREATE OR REPLACE TABLE raw_links (  
    movieId INTEGER,  
    imdbId INTEGER,  
    tmdbId INTEGER  
);  
COPY INTO raw_links FROM '@netflixstage/links.csv' FILE_FORMAT = (TYPE =  
'CSV' SKIP_HEADER = 1 FIELD_OPTIONALLY_ENCLOSED_BY = '');
```

✅ Validate row counts for each table to confirm ingestion.

3 Setting Up dbt Locally

Install dbt using pipx:

```
pipx install dbt-snowflake  
pipx ensurepath
```

Add ~/.local/bin to your PATH in ~/.zshrc:

```
export PATH="$PATH:/Users/your_username/.local/bin"  
source ~/.zshrc
```

Initialize the dbt Project

```
dbt init my_dbt_project
```

Configure ~/.dbt/profiles.yml:

```
my_dbt_project:  
  target: dev  
  outputs:  
    dev:  
      type: snowflake
```

```
account: your_account
user: dbt
password: your_password
role: TRANSFORM
database: MOVIELENS
warehouse: COMPUTE_WH
schema: RAW
threads: 1
```

✓ Validate with:

```
dbt debug
```

4 Create `dbt_project.yml`

```
# Name your project! Project names should contain only lowercase characters
# and underscores. A good package name should reflect your organization's
# name or the intended use of these models
name: 'my_dbt_project'
version: '1.0.0'

# This setting configures which "profile" dbt uses for this project.
profile: 'my_dbt_project'

# These configurations specify where dbt should look for different types of files.
# The `model-paths` config, for example, states that models in this project can be
# found in the "models/" directory. You probably won't need to change these!
model-paths: ["models"]
analysis-paths: ["analyses"]
test-paths: ["tests"]
seed-paths: ["seeds"]
macro-paths: ["macros"]
snapshot-paths: ["snapshots"]

clean-targets:          # directories to be removed by `dbt clean`
  - "target"
  - "dbt_packages"

# Configuring models
# Full documentation: https://docs.getdbt.com/docs/configuring-models

# In this example config, we tell dbt to build all models in the example/
# directory as views. These settings can be overridden in the individual model
# files using the `{{ config(...) }}` macro.
models:
```

```
my_dbt_project:
  +materialized: view
  dim:
    +materialized: table
  fct:
    +materialized: table
```

5 Building the dbt Pipeline

Create and populate models/staging/ with src_ models (views):

- src_movies.sql, src_ratings.sql, src_tags.sql, src_genome_tags.sql, src_genome_scores.sql, src_links.sql referencing MOVIELENS.RAW tables.

```
{% config(materialized='view') %}
•
• WITH raw_movies AS (
•   SELECT * FROM movielens.raw.raw_movies
• )
• SELECT
•   movieId AS movie_id,
•   title,
•   genres
• FROM raw_moviesdbt
```

```
WITH raw_ratings AS (
  SELECT * FROM MOVIELENS.RAW.raw_ratings
)

SELECT
  userId AS user_id,
  movieId AS movie_id,
  rating,
  TO_TIMESTAMP_LTZ(timestamp) AS rating_timestamp
FROM raw_ratings
```

```
WITH raw_tags AS (
  SELECT * FROM MOVIELENS.RAW.raw_tags
)

SELECT
  userId AS user_id,
```

```
movieId AS movie_id,  
tag,  
TO_TIMESTAMP_LTZ(timestamp) AS tag_timestamp  
FROM raw_tags
```

```
WITH raw_links AS (  
    SELECT * FROM MOVIELENS.RAW.raw_links  
)  
  
SELECT  
    movieId AS movie_id,  
    imdbId AS imdb_id,  
    tmdbId AS tmdb_id  
FROM raw_links
```

```
WITH raw_genome_tags AS (  
    SELECT * FROM MOVIELENS.RAW.raw_genome_tags  
)  
  
SELECT  
    tagId AS tag_id,  
    tag  
FROM raw_genome_tags
```



Run:

```
dbt run
```



Validate creation of SRC_* views.

Create and populate models/dim/ with dim_ models (tables):

Example dim_movies.sql:

```
{{ config(materialized='table') }}
```

```
WITH src_movies AS (
```



```

        SELECT * FROM {{ ref('src_movies') }}
    )
SELECT
    movie_id,
    INITCAP(TRIM(title)) AS movie_title,
    SPLIT(genres, '|') AS genre_array,
    genres
FROM src_movies

```

✅ Repeat similarly

for dim_ratings.sql, dim_tags.sql, dim_links.sql, dim_genome_tags.sql, dim_genome_scores.sql using {{ ref('src_*) }}

✅ Run:

```
dbt run
```

✅ Validate creation of DIM_* tables.

Create and populate models/fct/ with fct_models (fact tables):

Example fct_user_activity.sql:

```

{{ config(materialized='table') }}

WITH ratings AS (
    SELECT DISTINCT user_id FROM {{ ref('dim_ratings') }}
),
tags AS (
    SELECT DISTINCT user_id FROM {{ ref('dim_tags') }}
)
SELECT DISTINCT user_id
FROM (
    SELECT * FROM ratings
    UNION
    SELECT * FROM tags
)
)

```

✅ Run:

```
dbt run
```

✅ Validate creation of FCT_* fact tables.

6 Testing and Documentation

✓ Create tests in tests/ for not_null and unique constraints on primary keys.

✓ Run:

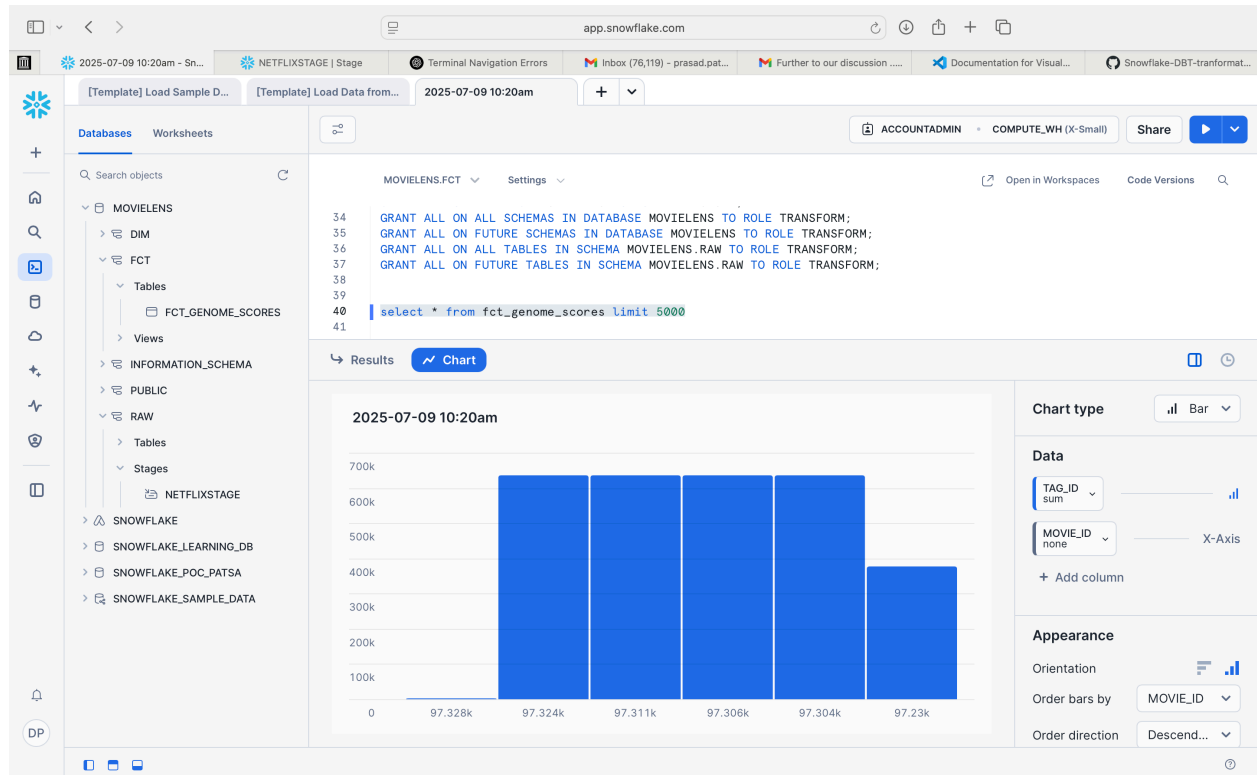
```
dbt test
```

The screenshot shows the Snowflake web interface. The left sidebar displays a tree view of databases and schemas, including MOVIELENS, DIM, FCT, INFORMATION_SCHEMA, PUBLIC, RAW, and STAGES. The main panel shows a SQL query in the MOVIELENS.FCT schema. The query is as follows:

```
34 GRANT ALL ON ALL SCHEMAS IN DATABASE MOVIELENS TO ROLE TRANSFORM;  
35 GRANT ALL ON FUTURE SCHEMAS IN DATABASE MOVIELENS TO ROLE TRANSFORM;  
36 GRANT ALL ON ALL TABLES IN SCHEMA MOVIELENS.RAW TO ROLE TRANSFORM;  
37 GRANT ALL ON FUTURE TABLES IN SCHEMA MOVIELENS.RAW TO ROLE TRANSFORM;  
38  
39  
40 select * from fct_genome_scores limit 5000  
41
```

Below the query, the 'Results' tab is selected, showing a table with 14 rows and 3 columns: MOVIE_ID, TAG_ID, and RELEVANCE_SCORE. The table data is as follows:

	# MOVIE_ID	# TAG_ID	# RELEVANCE_SCORE
1	97230	721	0.0608
2	97230	722	0.0195
3	97230	723	0.0112
4	97230	724	0.013
5	97230	725	0.0408
6	97230	726	0.1915
7	97230	727	0.1023
8	97230	728	0.029
9	97230	729	0.0235
10	97230	730	0.0782
11	97230	731	0.0448
12	97230	732	0.028
13	97230	733	0.7862
14	97230	734	0.035



✓ Generate documentation and explore lineage:

```
dbt docs generate
dbt docs serve
```

Catalog written to
/Users/**durgapatsa**/Documents/my_dbt_project/target/catalog.json
durgapatsa@MacBookPro my_dbt_project %

```
{"metadata": {"dbt_schema_version": "https://schemas.getdbt.com/dbt/catalog/v1.json",
"dbt_version": "1.10.3", "generated_at": "2025-07-09T22:32:57.109209Z",
"invocation_id": "64ac51d6-f813-4ef0-80af-20ef36aade1c", "invocation_started_at":
"2025-07-09T22:32:35.173983Z", "env": {}}, "nodes": {"model.my_dbt_project.src_tags":
{"metadata": {"type": "VIEW", "schema": "FCT", "name": "SRC_TAGS", "database":
"MOVIELENS", "comment": null, "owner": "TRANSFORM"}, "columns": {"USER_ID": {"type":
"NUMBER", "index": 1, "name": "USER_ID", "comment": null}, "MOVIE_ID": {"type":
"NUMBER", "index": 2, "name": "MOVIE_ID", "comment": null}, "TAG": {"type": "TEXT",
"index": 3, "name": "TAG", "comment": null}, "TAG_TIMESTAMP": {"type":
"TIMESTAMP_LTZ", "index": 4, "name": "TAG_TIMESTAMP", "comment": null}}, "stats":
{"has_stats": {"id": "has_stats", "label": "Has Stats?", "value": false, "include":
false, "description": "Indicates whether there are statistics for this table"}},
"unique_id": "model.my_dbt_project.src_tags"},
"model.my_dbt_project.fct_genome_scores": {"metadata": {"type": "BASE TABLE",
"schema": "FCT", "name": "FCT_GENOME_SCORES", "database": "MOVIELENS", "comment":
null, "owner": "TRANSFORM"}, "columns": {"MOVIE_ID": {"type": "NUMBER", "index": 1,
```

```
"name": "MOVIE_ID", "comment": null}, {"TAG_ID": {"type": "NUMBER", "index": 2, "name": "TAG_ID", "comment": null}, {"RELEVANCE_SCORE": {"type": "FLOAT", "index": 3, "name": "RELEVANCE_SCORE", "comment": null}}, {"stats": {"last_modified": {"id": "last_modified", "label": "Last Modified", "value": "2025-07-09 21:34UTC", "include": true, "description": "The timestamp for last update/change"}, {"bytes": {"id": "bytes", "label": "Approximate Size", "value": 46086144, "include": true, "description": "Approximate size of the table as reported by Snowflake"}, {"row_count": {"id": "row_count", "label": "Row Count", "value": 11709768, "include": true, "description": "An approximate count of rows in this table"}, {"has_stats": {"id": "has_stats", "label": "Has Stats?", "value": true, "include": false, "description": "Indicates whether there are statistics for this table"}}, {"unique_id": "model.my_dbt_project.fct_genome_scores"}}, {"model.my_dbt_project.src_genome_score": {"metadata": {"type": "VIEW", "schema": "FCT", "name": "SRC_GENOME_SCORE", "database": "MOVIELENS", "comment": null, "owner": "TRANSFORM"}, {"columns": {"MOVIE_ID": {"type": "NUMBER", "index": 1, "name": "MOVIE_ID", "comment": null}, {"TAG_ID": {"type": "NUMBER", "index": 2, "name": "TAG_ID", "comment": null}, {"RELEVANCE": {"type": "FLOAT", "index": 3, "name": "RELEVANCE", "comment": null}}, {"stats": {"has_stats": {"id": "has_stats", "label": "Has Stats?", "value": false, "include": false, "description": "Indicates whether there are statistics for this table"}}, {"unique_id": "model.my_dbt_project.src_genome_score"}}, {"model.my_dbt_project.src_movies": {"metadata": {"type": "VIEW", "schema": "FCT", "name": "SRC_MOVIES", "database": "MOVIELENS", "comment": null, "owner": "TRANSFORM"}, {"columns": {"MOVIE_ID": {"type": "NUMBER", "index": 1, "name": "MOVIE_ID", "comment": null}, {"TITLE": {"type": "TEXT", "index": 2, "name": "TITLE", "comment": null}, {"GENERES": {"type": "TEXT", "index": 3, "name": "GENERES", "comment": null}}, {"stats": {"has_stats": {"id": "has_stats", "label": "Has Stats?", "value": false, "include": false, "description": "Indicates whether there are statistics for this table"}}, {"unique_id": "model.my_dbt_project.src_movies"}}, {"model.my_dbt_project.src_links": {"metadata": {"type": "VIEW", "schema": "FCT", "name": "SRC_LINKS", "database": "MOVIELENS", "comment": null, "owner": "TRANSFORM"}, {"columns": {"MOVIE_ID": {"type": "NUMBER", "index": 1, "name": "MOVIE_ID", "comment": null}, {"IMDB_ID": {"type": "NUMBER", "index": 2, "name": "IMDB_ID", "comment": null}, {"TMDB_ID": {"type": "NUMBER", "index": 3, "name": "TMDB_ID", "comment": null}}, {"stats": {"has_stats": {"id": "has_stats", "label": "Has Stats?", "value": false, "include": false, "description": "Indicates whether there are statistics for this table"}}, {"unique_id": "model.my_dbt_project.src_links"}}, {"model.my_dbt_project.src_ratings": {"metadata": {"type": "VIEW", "schema": "FCT", "name": "SRC_RATINGS", "database": "MOVIELENS", "comment": null, "owner": "TRANSFORM"}, {"columns": {"USER_ID": {"type": "NUMBER", "index": 1, "name": "USER_ID", "comment": null}, {"MOVIE_ID": {"type": "NUMBER", "index": 2, "name": "MOVIE_ID", "comment": null}, {"RATING": {"type": "FLOAT", "index": 3, "name": "RATING", "comment": null}, {"RATING_TIMESTAMP": {"type": "TIMESTAMP_LTZ", "index": 4, "name": "RATING_TIMESTAMP", "comment": null}}, {"stats": {"has_stats": {"id": "has_stats", "label": "Has Stats?", "value": false, "include": false, "description": "Indicates whether there are statistics for this table"}}, {"unique_id": "model.my_dbt_project.src_ratings"}}, {"model.my_dbt_project.src_genome_tags": {"metadata": {"type": "VIEW", "schema": "FCT", "name": "SRC_GENOME_TAGS", "database": "MOVIELENS", "comment": null, "owner": "TRANSFORM"}, {"columns": {"TAG_ID": {"type": "NUMBER", "index": 1, "name": "TAG_ID", "comment": null}, {"TAG": {"type": "TEXT",
```

```
"index": 2, "name": "TAG", "comment": null}}, "stats": {"has_stats": {"id":  
"has_stats", "label": "Has Stats?", "value": false, "include": false, "description":  
"Indicates whether there are statistics for this table"}}, "unique_id":  
"model.my_dbt_project.src_genome_tags"}}, "sources": {}, "errors": null}
```

✅ Final Verification Checklist:

- ✅ Snowflake configured with raw data ingested from S3.
- ✅ dbt environment installed, debugged, and connected.
- ✅ staging models (SRC_*) created as views.
- ✅ dim models (DIM_*) created as tables.
- ✅ fct models (FCT_*) created as tables.
- ✅ Tests passing.
- ✅ Lineage documented and visualized.
- ✅ Ready for ML and dashboard integrations.



Next Steps

- Create Power BI / Streamlit dashboards using FCT_* tables.
- Integrate GitHub Actions for CI/CD on your dbt project.
- Extend with incremental models for large-scale data pipelines.

This structured end-to-end pipeline prepares you to **showcase your Snowflake + dbt skills confidently for your Data Architect or Senior Data Engineer roles** while reinforcing systematic architecture practices for your real-world learning.