# Data Integration and Transformation Pipeline

# Content

1 Into. &
Task Management

2 Snowflake
Environmental
Setup

3 Snowflake
Data Pipeline
Architecture

4 Dashboard for
Analysis

5 DBT
Transformation

6 Retrospective

# Brainstorming

# Snowflake Environmental Setup

**1**

Create Team
Warehouse

**Create Database
and Schemas**

**2**

· RAW
· PREP
· REFINED
· DELIVERY

**3**

**Grant Privileges**

# Grant Privileges

## Privileges

---

👤 CHIPMUNK_ROLE  (Current Role)  `USAGE`

---

👤 KOALA_ROLE  `USAGE`

---

👤 LEMMING_ROLE  `USAGE`

---

👤 LEMUR_ROLE  `USAGE`

---

👤 TEAM_4_USER_ROLE  `🔍 OWNERSHIP`  `DELETE - FUTURE TABLE`  `INSERT - FUTURE TABLE`  `SELECT - FUTURE TABLE`  `SELECT - FUTURE VIEW`  `UPDATE - FUTURE TABLE`  `USAGE - FUTURE SCHEMA`

---

👤 TEAM_4_VIEWER_ROLE  `SELECT - FUTURE VIEW`  `USAGE`

# Snowflake Data Pipeline Architecture

RAW ➜ PREP ➜ REFINED ➜ DELIVERY

CSV

JSON

Marketplace

Extract    Load

RAW Data

**Data Warehouse**

Transform

Transform

Transform

Analyze

# Extract



**Housing rental information data**

Kaggle

**Quality of Life Index Metrics**

Kaggle

**American Community Survey**

Marketplace

**Load**    CSV / Json →    **TEAM_4_DB.RAW**

**CSV**
Housing rental information data

**JSON**
Quality of Life Index Metrics

American Community Survey

SQL → CSV

```
C:\Users\sugan>snowsql -a nlb11398 -u KOALA -r KOALA_ROLE -d TEAM_4_DB -s RAW
* SnowSQL * v1.4.1
Type SQL statements or !help
KOALA#LEARNER_WH@TEAM_4_DB.RAW>USE WAREHOUSE TEAM_4_WAREHOUSE;
+-----------------------------------+
| status                            |
|-----------------------------------|
| Statement executed successfully.  |
+-----------------------------------+
1 Row(s) produced. Time Elapsed: 0.246s
KOALA#TEAM_4_WAREHOUSE@TEAM_4_DB.RAW>SELECT CURRENT_WAREHOUSE();
+---------------------+
| CURRENT_WAREHOUSE() |
|---------------------|
| TEAM_4_WAREHOUSE    |
+---------------------+
1 Row(s) produced. Time Elapsed: 0.233s
KOALA#TEAM_4_WAREHOUSE@TEAM_4_DB.RAW>PUT
file://C:/Users/sugan/OneDrive/Desktop/Hyper_Island/Course_8_Data_Engineering/
```
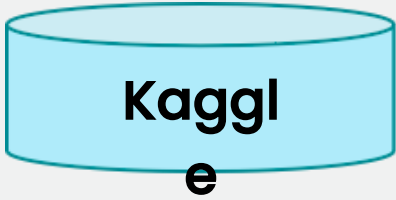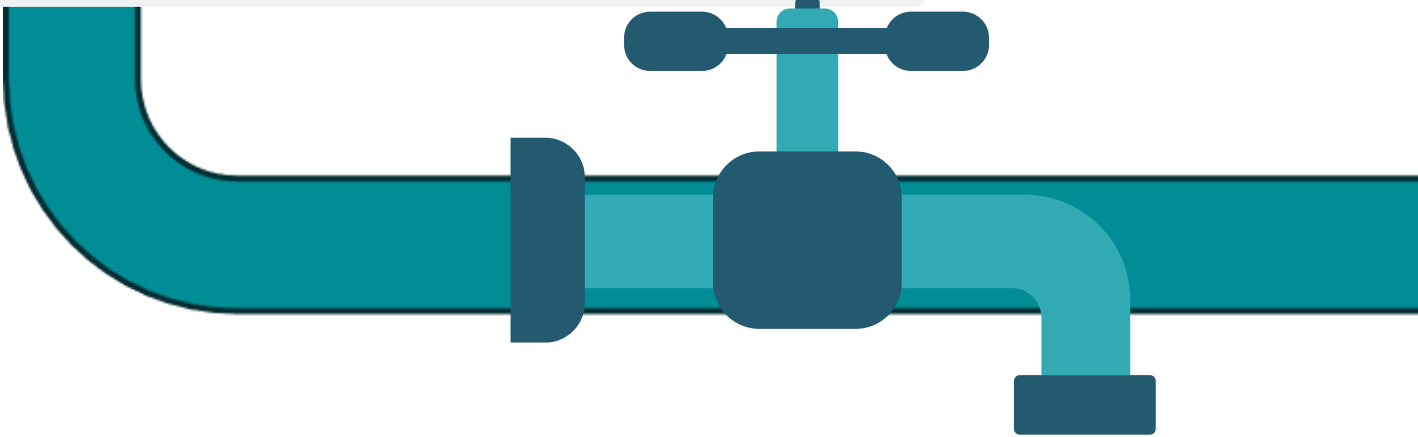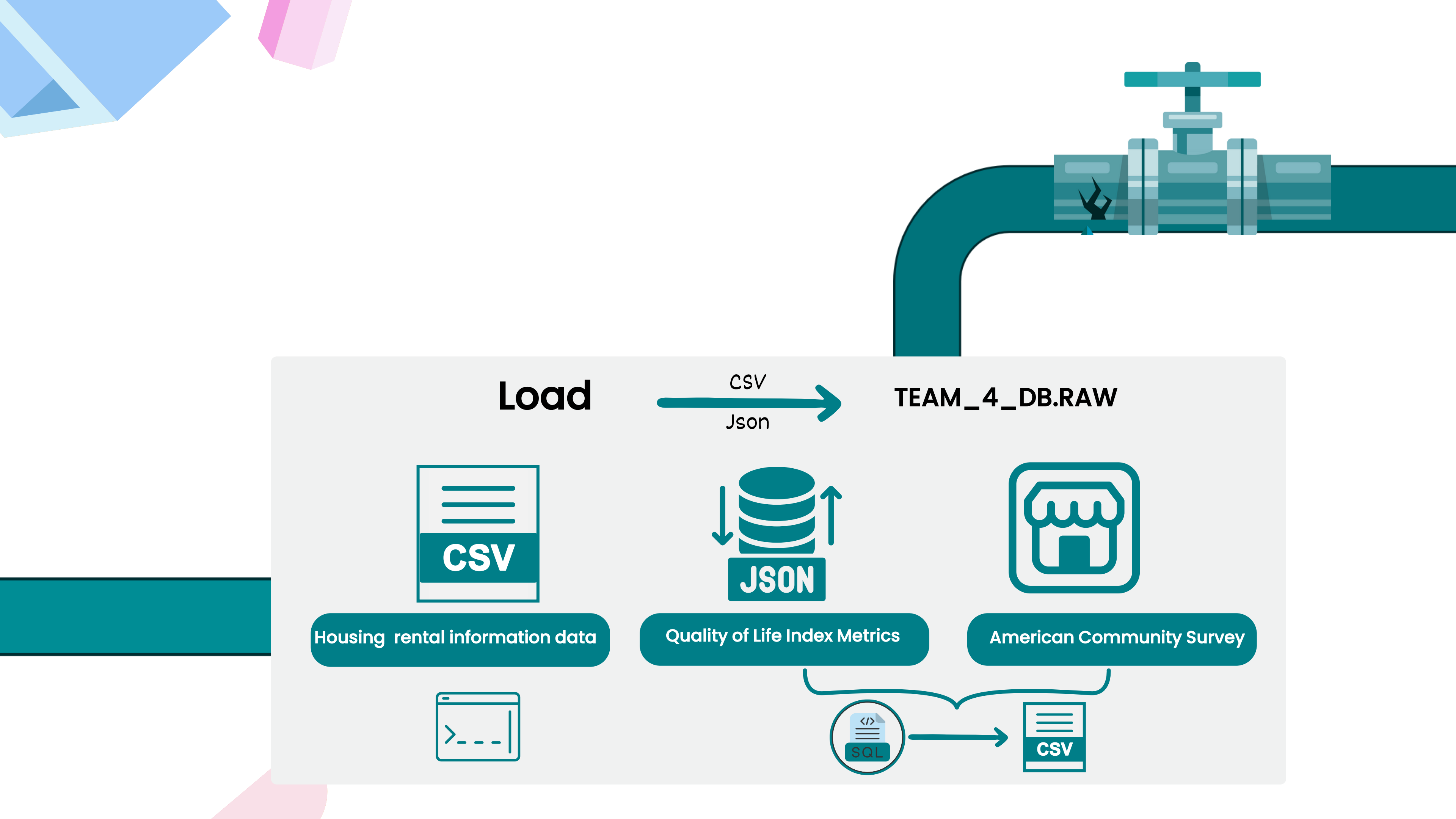
#Raw

## Transform

- Standardized column names
- Limited to the needed columns
- Type conversion
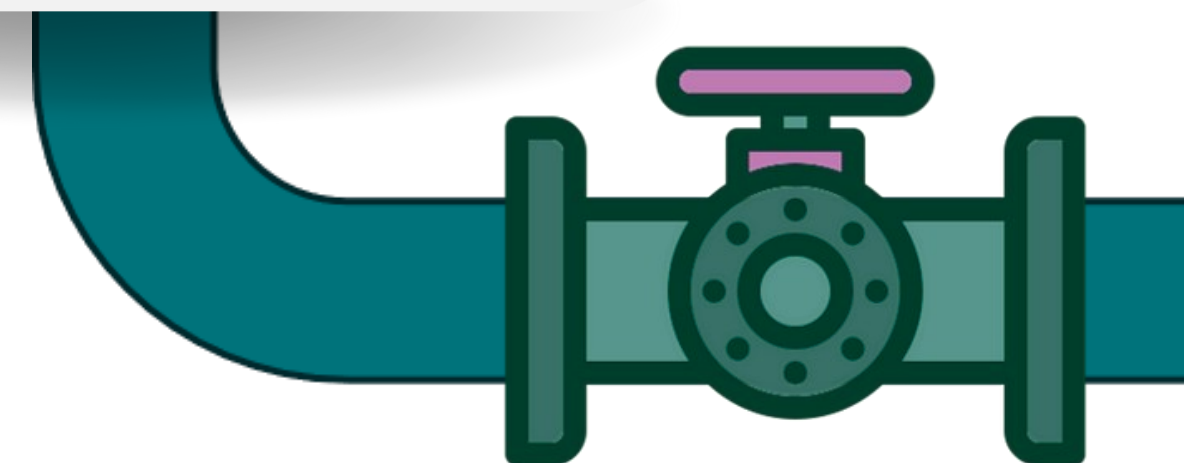- Parsed JSON to structured columns

#Raw ➡️ Prep

## Transform

- Data Cleaning:
- Standardizing State Codes
- Removing outliers
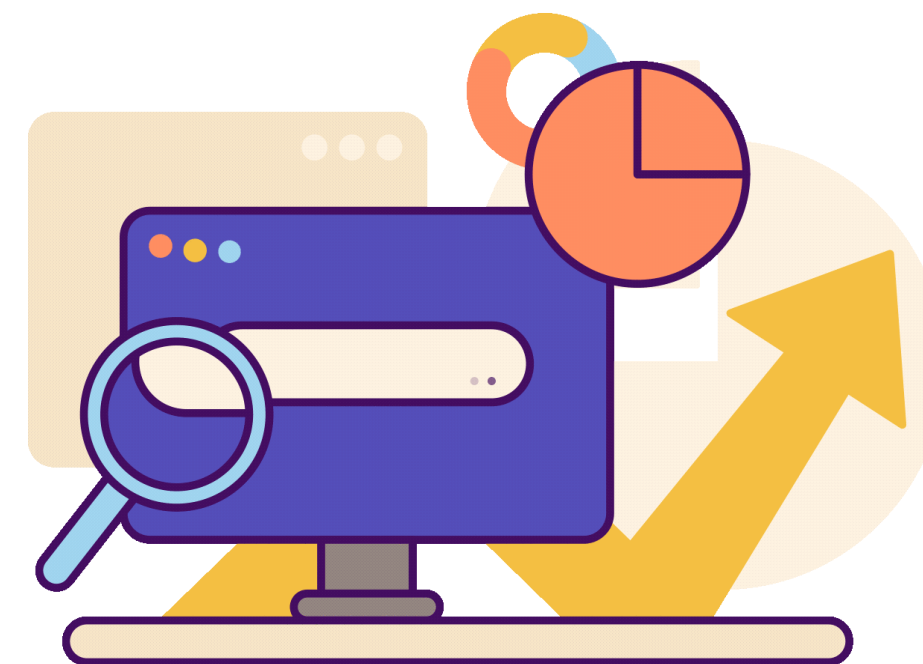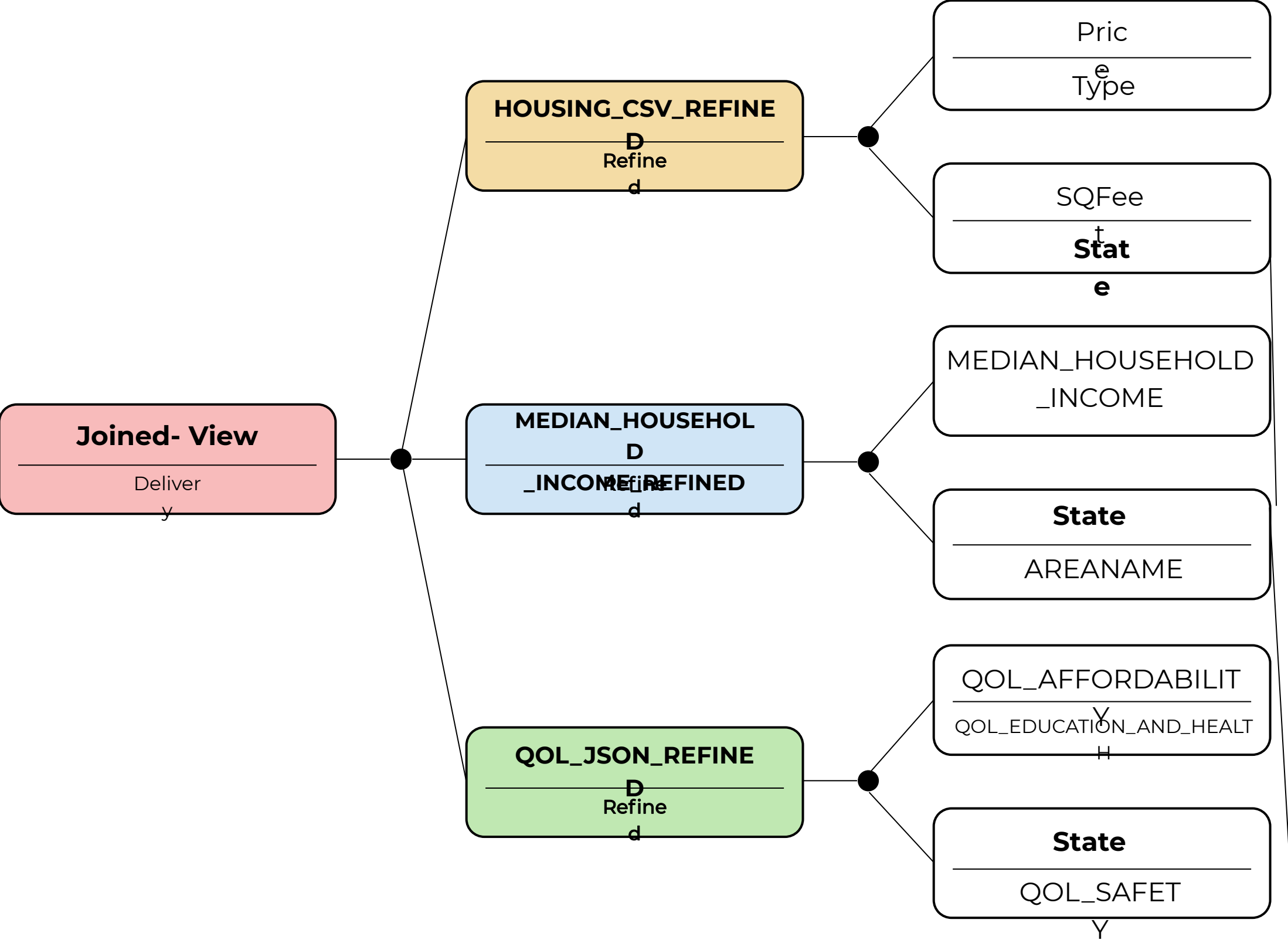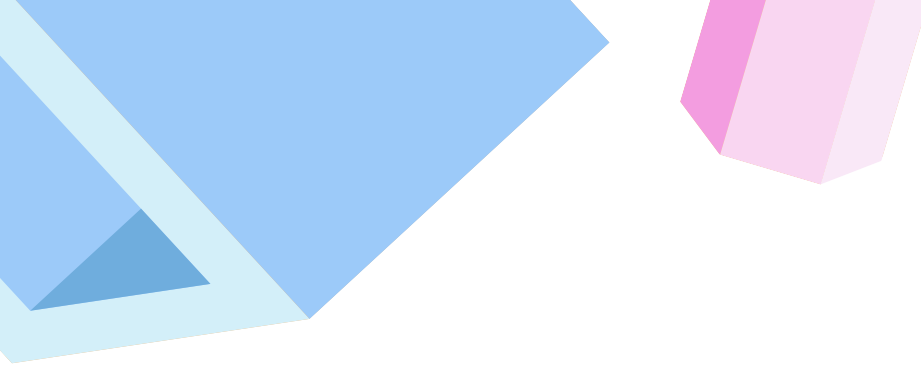
#Raw ➡️ Prep ➡️ Refined

## Dashboard ❄️

- Using Team_4_Viewer_role to create dashboard
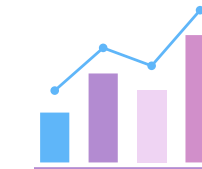- From refined schema joined 3 tables together for analysis.

#Raw ➡️ Prep ➡️ Refined ➡️ Delivery

# Snowflake Dashboard
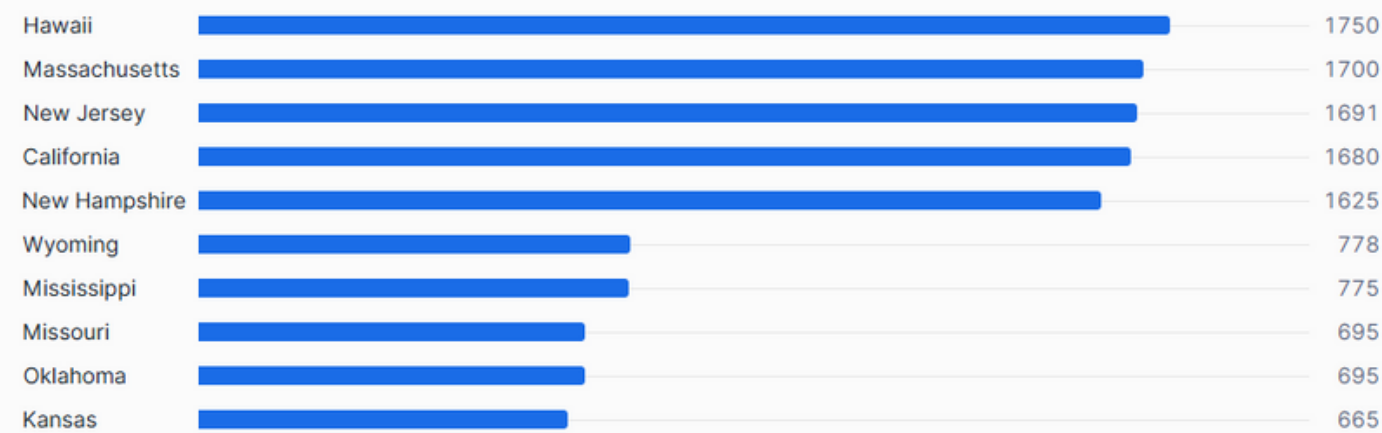
TEAM_4_VIEWER_ROLE • TEAM_4_WAREHOUSE (X-Small)   Share   ▶ Run

+ ⚏

Updated 11h ago

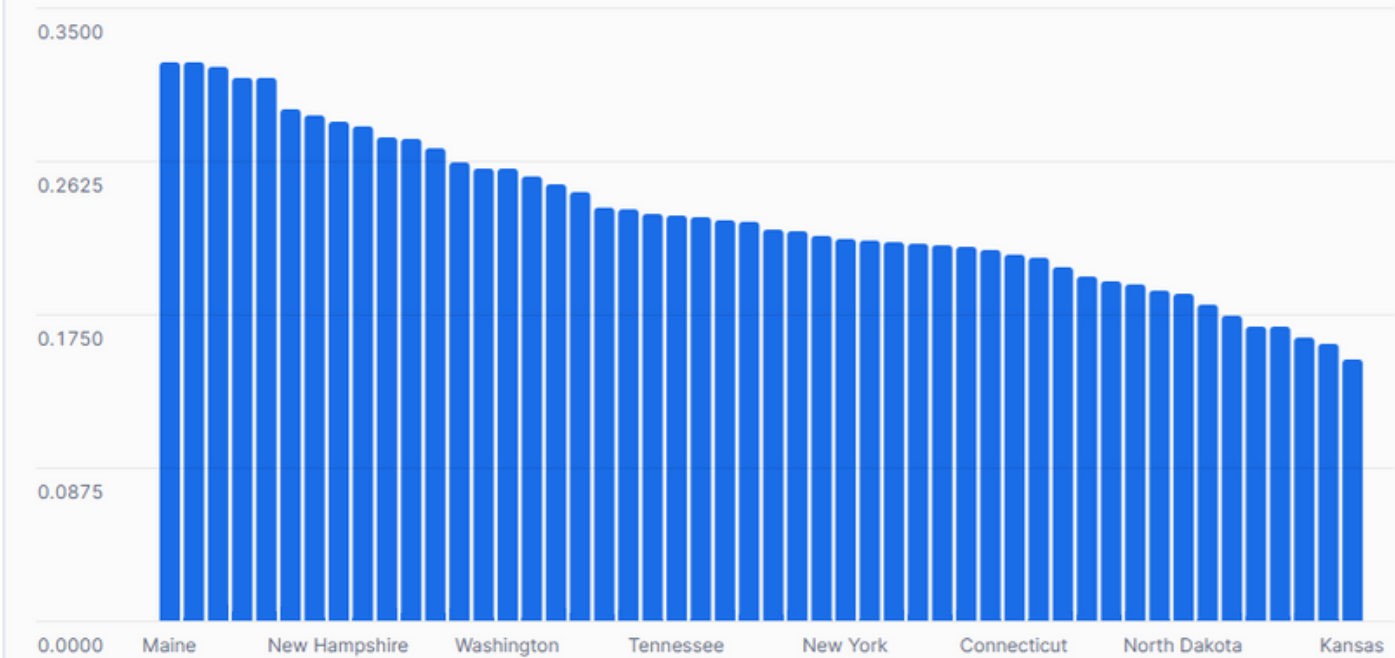## Top 5 & Bottom 5 States By Monthly Rent Price  ⋯

| State | Rent |
|---|---|
| Hawaii | 1750 |
| Massachusetts | 1700 |
| New Jersey | 1691 |
| California | 1680 |
| New Hampshire | 1625 |
| Wyoming | 778 |
| Mississippi | 775 |
| Missouri | 695 |
| Oklahoma | 695 |
| Kansas | 665 |

## Rent_Price_Influencer  ⋯

0.8
0.4
0.0
-0.4
-0.8

Income-Rent   Education-Health-Rent   Economy-Rent   Safety-Rent   Affordability-Rent

## Average Rent by Size Group & Type  ⋯

● apartment  ● condo  ● duplex

2000
1500
1000
500.0
0.0

Large (>1000 sqft)   Medium (500–1000 sqft)   Small (<500 sqft)

## Rend Burden  ⋯

0.3500
0.2625
0.1750
0.0875
0.0000

Maine   New Hampshire   Washington   Tennessee   New York   Connecticut   North Dakota   Kansas

# Integrated pipeline

# dbt Snowflake Pipeline Architecture

- Transformation
- Data Tests

**dbt**

**PREP**
**REFINED**

CSV

JSON

Marketplace

**RAW**

**DELIVERY**

Snowflake

Raw Data

Transformed Data

dbt

# GitHub collaboration

## Change Branch

Git Branch

main — updated Mon Jun 09 2025 ⌄

main — updated Mon Jun 09 2025 ✓

suganyam2001-housing — updated Mon Jun 09 2025 (current branch)

## Users

🔍 Search users by name or email

| Name ↑ | License |
|---|---|
| **Danqing Yao**<br>danqing.yao@hyperisland.se | Developer |
| **Suganya Muruganantham**<br>suganya.m2001@gmail.com | Developer |
| **Sunny Gustavsson**<br>sunny.zhang.qing@gmail.com | Developer |
| **Tanglan Yang**<br>volingla@gmail.com | Developer |

# DAG

# Retrospective

# RBAC (Role-Based Access Control) Implementation

**Initial Setup:** Created database, schemas, and warehouse using Training Role.

Challenge: Training Role had broad privileges.

Team members could grant themselves more access even after restrictions.

Solution: Created team-specific role for collaboration.

Lesson: Aligned with least privilege for secure access.

# RBAC (Role-Based Access Control) Implementation

# Building an End-to-End Data Pipeline in Snowflake and DBT Integrated

**What worked well:**
Learning the entire pipeline in Snowflake and DBT

Challenge: Redundant transformation

Understanding Integration took time

Key Learning: Modern best practice and DBT with Git ensures automation control

Lesson for Next time: Plan, start small and scale