

# US Demographics and Public Resources

---

Sashimi

# Chosen Datasets – Primary

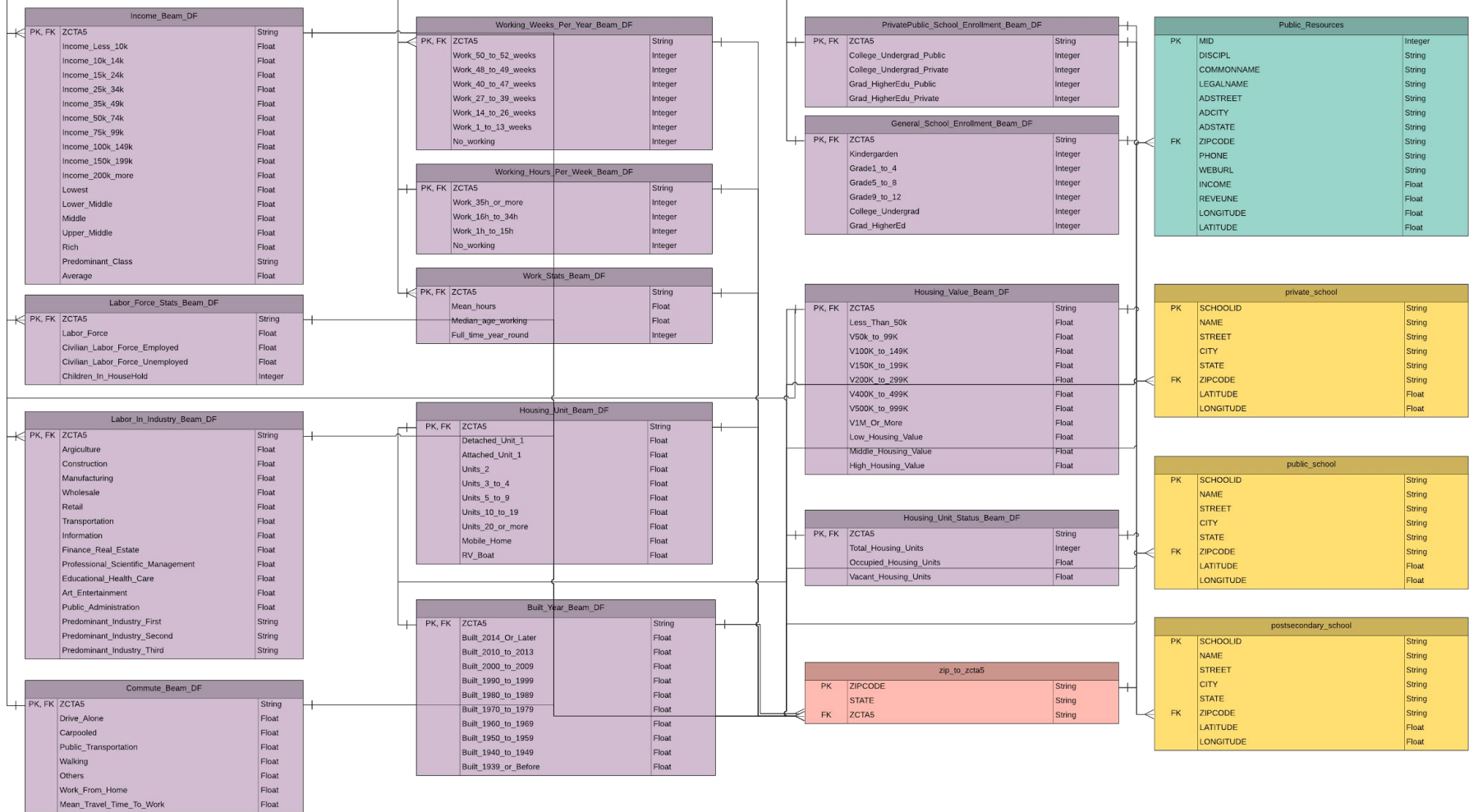
---

- American Housing Survey
  - Basic demographic, household income
  - Housing status and information
  - Predominant industry, employment rate, average working hour/week
  - Overall education level, school enrollment rate
  - Location in ZCTA5

# Chosen Datasets – Secondary

---

- Institute of Library and Museums
  - Discipline, income, revenue
  - Location in ZIPCODE, longitude, latitude
- Education Demographic and Geographic Estimates
  - Public, private and postsecondary school
  - Location in ZIPCODE, longitude, latitude
- UDS Mapper
  - Location in ZIPCODE
  - Location in ZCTA5



# Area of Interest

---

- Average housing value vs average household income
- Industry vs average household income
- Which area has better public resources vs which area is lacking
- Investment vs demand on education resources
- Where public resources located at vs living condition nearby

# Problems and Solutions

---

- Massive Table
- Column names
- Different location keys (Zip Code vs. ZCTA5)
- Omit unnecessary columns and split up table for normalization
- Beam transform
- UDS Mapper to connect

# Pipelines

---

- SQL
  - Cast special null markers as 'Null'
  - Update schema type
- Beam
  - Update column names from code to human readable name
  - Generate columns through computation
- DAG

# Live Demo

---

- Income: generating average and assigning socioeconomic classes
- Cross joins
- Visualization



# Future Improvements

---

- More comprehensive “public resources”
  - Combine Schools with Libraries and Museums
- Collect some data on government funding
- Zip Code might be too granular when aggregating data
  - Example: I live in Riverside but go to UT
- Better estimate on the needed data to avoid unnecessary work

Thanks & Questions