

# UWA DATA ANALYSIS BOOTCAMP - PROJECT 2 ETL

## TEAM 6 PROPOSAL

### Our Team



Johan Snyder

Johan Snyder  
JohanfromEsperance



Jon Wood

Jono Wood  
jonowood

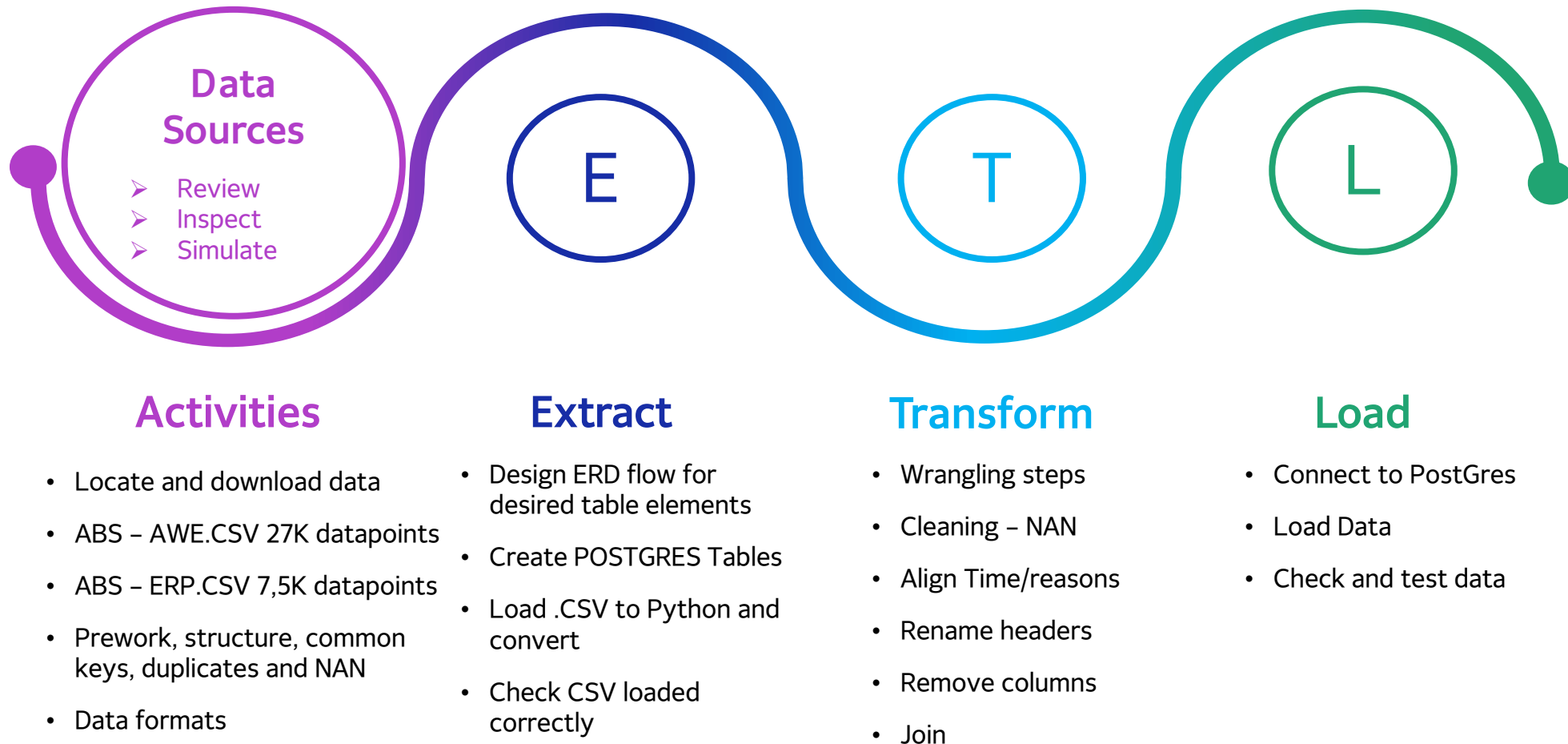
### Our Git Repository

[https://github.com/jonowood/Project\\_2](https://github.com/jonowood/Project_2)

# PROJECT 2-ETL : FINANCE/SOCIAL SERVICES

## HOW DOES AVERAGE WEEKLY EARNINGS INFLUENCE ESTIMATED RESIDENCE POPULATION

THE AIM OF OUR PROJECT IS TO UNCOVER PATTERNS IN POPULATION MOVEMENTS AND AVERAGE INCOME CHANGES. WE'LL EXAMINE RELATIONSHIPS BETWEEN POPULATION CHANGE AND AVERAGE INCOME, PERCENTAGE CHANGE AND OTHER RELATED RELATIONSHIPS DERIVED FROM THE DATA.



# EXTRACT

## Data Sets

Datasets Sourced From - <https://explore.data.abs.gov.au/> (Licence CC01)

File A - ABS - AWE.CSV (27K datapoints) - Population Movement Data for Australia

File B - ABS - ERP.CSV 7,5K (datapoints) - Average Income Data for Australia

## Initial Findings

Both data sets will require filtering of data based off multiple columns. The data is in a similar format across both sets.

Data is incomplete and will require the filling of NaN values. Some column values will require extra characters removed prior to transforming.

For loading this data to Postgres in a relational database, we will need to explode the data into multiple tables and link with date and location ID's.

# TRANSFORM

## REQUIRED STEPS TO TRANSFORM DATA;

File A - ABS - AWE.CSV (27K datapoints)

- Filter column C values to only include 'Earnings'
- Filter column B to only include 'All employees average weekly earnings'
- Filter column G to only include 'Original'
- Copy columns to DataFrame - REGION: Region, TIME\_PERIOD: Time Period, OBS\_VALUE

File B - ABS - ERP.CSV 7,5K (datapoints)

- Filter column B to only include '4: Internal Arrivals', '5: Internal Departures', '6: Net Internal Migration', '13: Change Over Previous Quarter'
- Copy columns to DataFrame - REGION: Region, TIME\_PERIOD: Time Period, OBS\_VALUE

Additional Pandas Transformations – Data Filtering, Data Mapping, Data Deduplication, Derived Variables

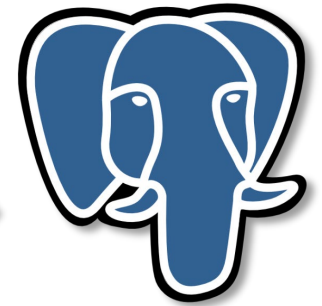
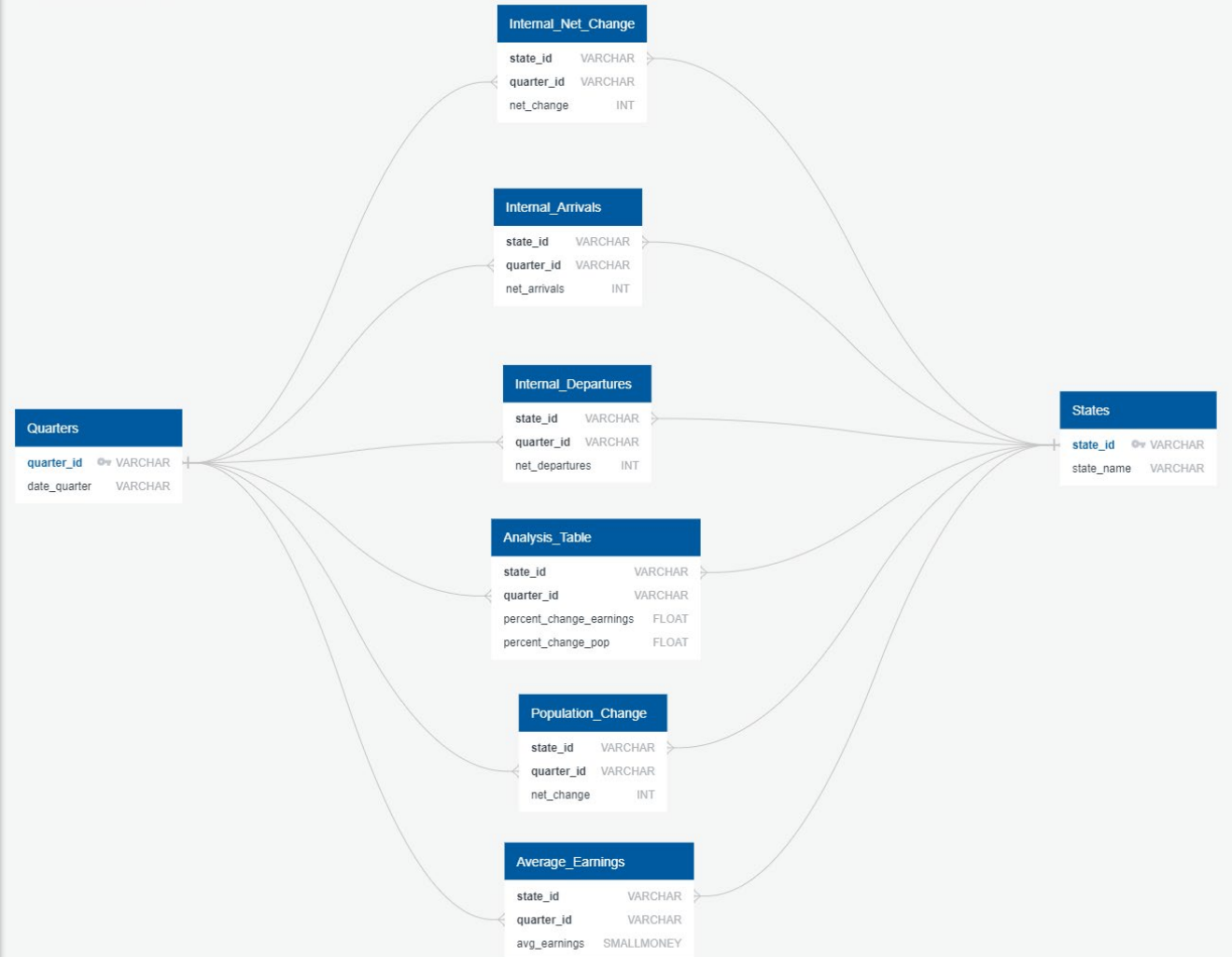
# LOAD

Relational Database - We will create a Database and Schema in PostgreSQL then Load the data for future Analysis

## Project 2

```
1
2
3 States
4 -----
5 state_id VARCHAR PK
6 state_name VARCHAR
7
8 Quarters
9 -----
10 quarter_id VARCHAR PK
11 date_quarter VARCHAR
12
13 Average_Earnings
14 -----
15 state_id VARCHAR FK >- States.state_id
16 quarter_id VARCHAR FK >- Quarters.quarter_id
17 avg_earnings SMALLMONEY
18
19 Population_Change
20 -----
21 state_id VARCHAR FK >- States.state_id
22 quarter_id VARCHAR FK >- Quarters.quarter_id
23 net_change INT
24
25 Internal_Arrivals
26 -----
27 state_id VARCHAR FK >- States.state_id
28 quarter_id VARCHAR FK >- Quarters.quarter_id
29 net_arrivals INT
30
31 Internal_Departures
32 -----
33 state_id VARCHAR FK >- States.state_id
34 quarter_id VARCHAR FK >- Quarters.quarter_id
35 net_departures INT
36
37 Internal_Net_Change
38 -----
39 state_id VARCHAR FK >- States.state_id
40 quarter_id VARCHAR FK >- Quarters.quarter_id
41 net_change INT
42
43 Analysis_Table
44 -----
45 state_id VARCHAR FK >- States.state_id
46 quarter_id VARCHAR FK >- Quarters.quarter_id
47 percent_change_earnings FLOAT
48 percent_change_pop FLOAT
49
50
```

www.quickdatabasediagrams.com



PostgreSQL