# ETL Project

## Real Estate/Walk Score Listings in Calgary

# Objective

- Extracting Real Estate data in Calgary, Alberta and Walk Scores for corresponding house addresses

- Walk Score: shows a walk/transit/bike score (0-100) for any address to the local downtown

- Transforming retrieved data into easy-to-read tables

- Loading transformed data into relational and non-relational databases for optimal functionality

# Data Sources

- Remax Canada: https://www.remax.ca/ab/calgary-real-estate
- Walk Score: https://www.walkscore.com/CA-AB/Calgary

# Extract

# Scraping Calgary Real Estate Data

**Remax**

```
In [1]: import pandas as pd
        import numpy as np
        import requests
        from bs4 import BeautifulSoup
        import time
        from splinter import Browser
        from sqlalchemy import create_engine
        import warnings
        warnings.filterwarnings('ignore')
        print('Libraries imported!')

        Libraries imported!
```

Using BeautifulSoup to scrape property details (house address, house details).

```
In [2]: house_address = []
        house_details = []

        base_url = 'https://www.remax.ca/ab/calgary-real-estate?page='
        urls = [base_url + str(x) for x in range(1,301)]
        time.sleep(2)
        for url in urls:
            # Parse HTML with Beautiful Soup
            time.sleep(2)
            response = requests.get(url)
            soup = BeautifulSoup(response.text, 'html.parser')

            try:
                addresses = soup.find_all('div', class_='left-content flex-one')
                for address in addresses:
                    house_address.append(address.text)
            except:
                house_address.append('None')

            try:
                details = soup.find_all('div', class_='property-details')
                for detail in details:
                    house_details.append(detail.text)
            except:
                house_details.append('None')
```

# Scraping Walk Score Data

```python
In [ ]:  scores_walk = []
         scores_bike = []
         scores_transit = []

         for i in post_code_list:

             try:
                 postal_code = i.replace(" ", "%20")
                 url_score = "https://www.walkscore.com/score/" + str(postal_code)
                 time.sleep(2)

                 # Parse HTML with Beautiful Soup
                 response = requests.get(url_score)
                 code_soup = BeautifulSoup(response.text, 'html.parser')

                 if 'pp.walk.sc/badge/walk/score' in str(code_soup):
                     ws = str(code_soup).split('pp.walk.sc/badge/walk/score/')[1][:2].replace('.','')
                     scores_walk.append(ws)
                 else:
                     ws = 'N/A'
                     scores_walk.append(ws)
                 if 'pp.walk.sc/badge/bike/score' in str(code_soup):
                     bs = str(code_soup).split('pp.walk.sc/badge/bike/score/')[1][:2].replace('.','')
                     scores_bike.append(bs)
                 else:
                     bs = 'N/A'
                     scores_bike.append(bs)
                 if 'pp.walk.sc/badge/transit/score' in str(code_soup):
                     ts = str(code_soup).split('pp.walk.sc/badge/transit/score/')[1][:2].replace('.','')
                     scores_transit.append(ts)
                 else:
                     ts = 'N/A'
                     scores_transit.append(ts)
             except:
                 ws = 'N/A'
                 scores_walk.append(ws)
                 bs = 'N/A'
                 scores_bike.append(bs)
                 ts = 'N/A'
                 scores_transit.append(ts)
```

# Transform

# Cleaning the Calgary Real Estate Data

## First dataframe: Address and Price details

```
In [3]: address_df = pd.DataFrame(house_address)

        new_df = address_df[0].str.split(' ', 2, expand=True)
        new_df["price"] = new_df[1].str.replace("$", "")
        new_df["price"] = new_df["price"].str.replace(",", "")
        new_df["price"] = pd.to_numeric(new_df["price"])

        del new_df[0]
        del new_df[1]
        new_df.head()
```

Out[3]:

|   | 2 | price |
|---|---|---|
| 0 | 9803 ELBOW DR SW, Calgary, AB, T2V 1M4 | 489900 |
| 1 | 101 - 3704 15A ST SW, Calgary, AB, T2T 4C3 | 319900 |
| 2 | 25 HARVEST GLEN WAY NE, Calgary, AB, T3K 4J2 | 399900 |
| 3 | 32 EVERGLEN GROVE SW, Calgary, AB, T2Y 4Z3 | 429500 |
| 4 | 416 THORNDALE RD NW, Calgary, AB, T2K 3C5 | 484900 |

# Cleaning the Calgary Real Estate Data

- Values separated into columns: price, address, postal code, bedrooms, bath, property type

- Price column type changed to integer

```
In [5]:  final_df = new_df[2].str.split(', Calgary, AB, ', expand=True)
         final_df.head()
```

Out[5]:

|   | 0 | 1 |
|---|---|---|
| 0 | 9803 ELBOW DR SW | T2V 1M4 |
| 1 | 101 - 3704 15A ST SW | T2T 4C3 |
| 2 | 25 HARVEST GLEN WAY NE | T3K 4J2 |
| 3 | 32 EVERGLEN GROVE SW | T2Y 4Z3 |
| 4 | 416 THORNDALE RD NW | T2K 3C5 |

```
In [6]:  df_add = pd.concat([new_df, final_df], axis=1)
         del df_add[2]
         df_add.columns = ["price", "address", "postal_code"]
         df_add.head()
```

Out[6]:

|   | price | address | postal_code |
|---|-------|---------|-------------|
| 0 | 489900 | 9803 ELBOW DR SW | T2V 1M4 |
| 1 | 319900 | 101 - 3704 15A ST SW | T2T 4C3 |
| 2 | 399900 | 25 HARVEST GLEN WAY NE | T3K 4J2 |
| 3 | 429500 | 32 EVERGLEN GROVE SW | T2Y 4Z3 |
| 4 | 484900 | 416 THORNDALE RD NW | T2K 3C5 |

# Cleaning the Calgary Real Estate Data

- Second dataframe: House details

```
In [7]:   details = pd.DataFrame(house_details)

          details_df = details[0].str.split('|', expand=True)
          details_df

          del details_df[2]

          details_df.columns = ["bedrooms", "bath", "property_type"]
          details_df.head()
```

Out[7]:

|   | bedrooms | bath | property_type |
|---|----------|------|---------------|
| 0 | 4 bed | 2 bath | house |
| 1 | 2 bed | 1 + 1 bath | condo |
| 2 | 3 bed | 2 bath | house |
| 3 | 3 bed | 2 + 1 bath | house |
| 4 | 3 bed | 2 bath | house |

# Joining the Calgary Real Estate Dataframes

- Concatenating House Address/Price details and House details dataframes.

```
In [8]:  calgary_df = pd.concat([df_add, details_df], axis=1)
         calgary_df.head()
```

Out[8]:

| | price | address | postal_code | bedrooms | bath | property_type |
|---|---|---|---|---|---|---|
| 0 | 489900 | 9803 ELBOW DR SW | T2V 1M4 | 4 bed | 2 bath | house |
| 1 | 319900 | 101 - 3704 15A ST SW | T2T 4C3 | 2 bed | 1 + 1 bath | condo |
| 2 | 399900 | 25 HARVEST GLEN WAY NE | T3K 4J2 | 3 bed | 2 bath | house |
| 3 | 429500 | 32 EVERGLEN GROVE SW | T2Y 4Z3 | 3 bed | 2 + 1 bath | house |
| 4 | 484900 | 416 THORNDALE RD NW | T2K 3C5 | 3 bed | 2 bath | house |

```
In [9]:  calgary_df.to_csv('calgary_df.csv', index=False)
```

# Cleaning the Walk Score Data

- Data converted into dataframe and columns named.

```
In [ ]:  score_df_trans = {'postal_code':postal_code_list, 'walk_score':scores_walk, 'bike_score':scores_bike, 'transit_score':s
         score_df_dup = pd.DataFrame(score_df_trans)
         score_df = score_df_dup.drop_duplicates()
         score_df.head()
```

```
In [19]: score_df.to_csv('score_df.csv', index=False)
```

```
In [9]:  score_df.head()
```

Out[9]:

|   | postal_code | walk_score | bike_score | transit_score |
|---|-------------|------------|------------|---------------|
| 0 | T2V 1M4 | 58.0 | 61.0 | 55.0 |
| 1 | T2T 4C3 | 53.0 | 81.0 | 42.0 |
| 2 | T3K 4J2 | 19.0 | 59.0 | 38.0 |
| 3 | T2Y 4Z3 | 6.0 | 30.0 | 31.0 |
| 4 | T2K 3C5 | 61.0 | 81.0 | 53.0 |

# Load
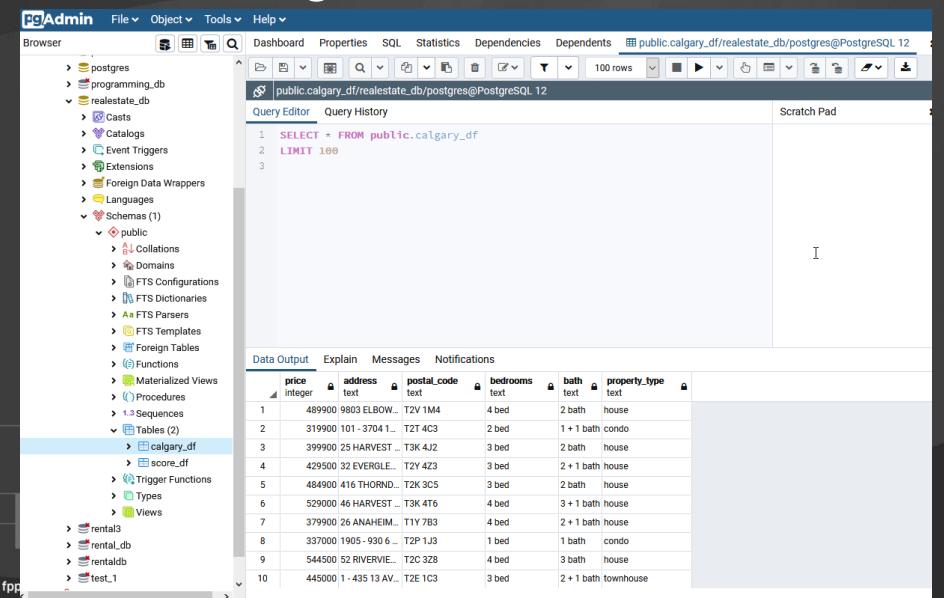
# Loading Data to Databases

- Building connection to PostgreSQL/MongoDB and loading transformed data
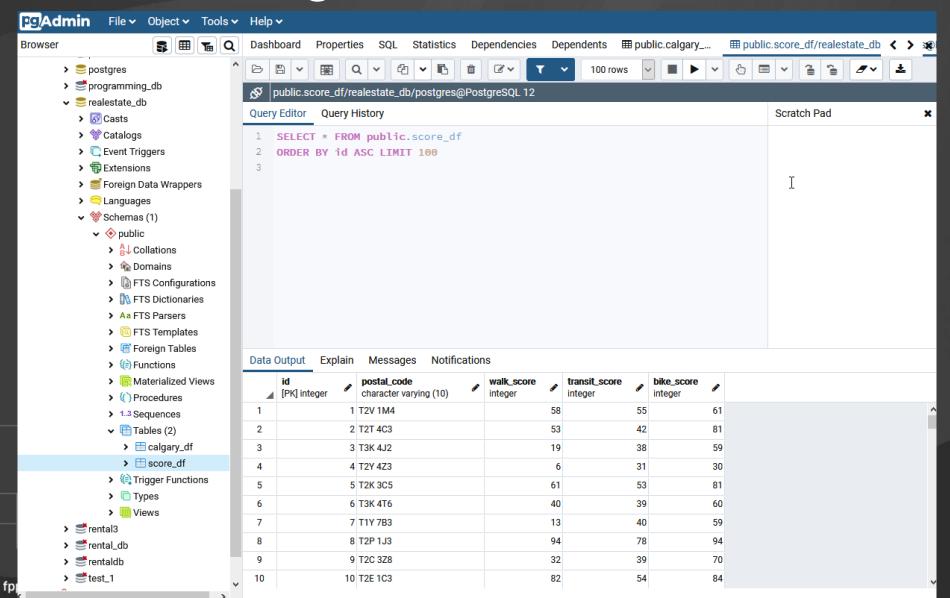
**SQL**

```
In [29]: calgary_df = pd.read_csv('calgary_df.csv')
         score_df = pd.read_csv('score_df.csv')

In [30]: rds_connection_string = "postgres:1@localhost:5432/realestate_db"
         engine = create_engine(f'postgresql://{rds_connection_string}')

         calgary_df.to_sql(name= "calgary_df", con=engine, if_exists="append", index=False)
         score_df.to_sql(name= "score_df", con=engine, if_exists="append", index=False)
```

# PostgreSQL Database

# PostgreSQL Database

# MongoDB Database

## MongoDB

```
In [ ]:  # Make a connection
         conn = "mongodb://localhost:27017"

         # Making a Connection with MongoClient
         client = MongoClient(conn)

         # database
         db = client.realestate_db


         collection = db.calgary
         calgary_dict = calgary_df.to_dict("records")
         collection.insert_many(calgary_dict)


         collection = db.score
         score_dict = score_df.to_dict("records")
         collection.insert_many(score_dict)
```
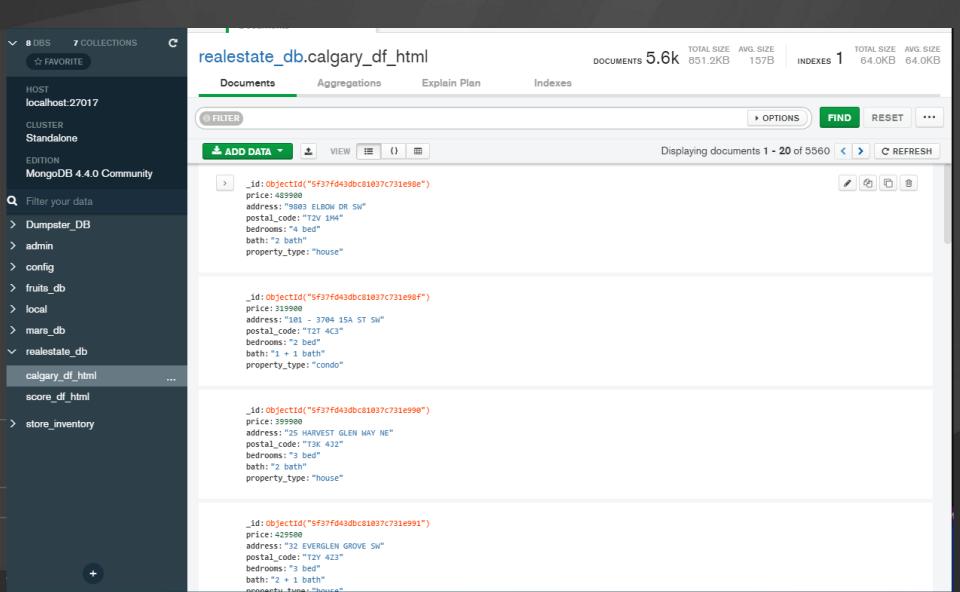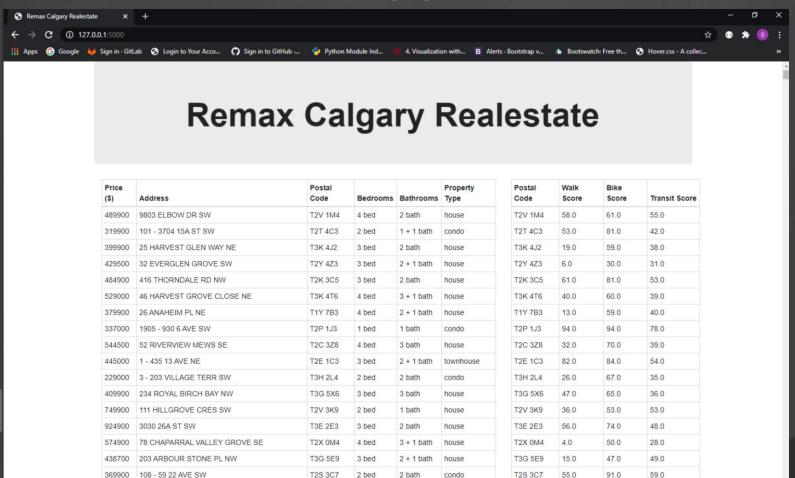
# MongoDB Database

# Converted Database into a Web Based Application