# National Housing Value Distribution

**Utilizing Zillow Data**

By: Gustavo Carmona, Kamila Czosnyka,
Sam Duffy, Caila Hanus

# Project Overview

- Engineered a datafile with over 1 million lines of data! (1,048,576 lines to be exact!)
- Used Zillow CSV files and API calls, we used ETL processes to clean, edit, and ultimately upload the data to SQL
- Used data to compare Chicago-specific housing values to major US cities
- Used data to compare Chicago-area housing data to itself
- This data could be used to look at historical changes in a certain area or to compare different regions

# API Resources / Data Extraction

- https://data.nasdaq.com/databases/ZILLOW
- This was used to both call data directly from the API as well as download pre-loaded CSV files
- The CSV files were needed as they provided information on the different Indicator and Region IDs which were needed in the URL of the API call.
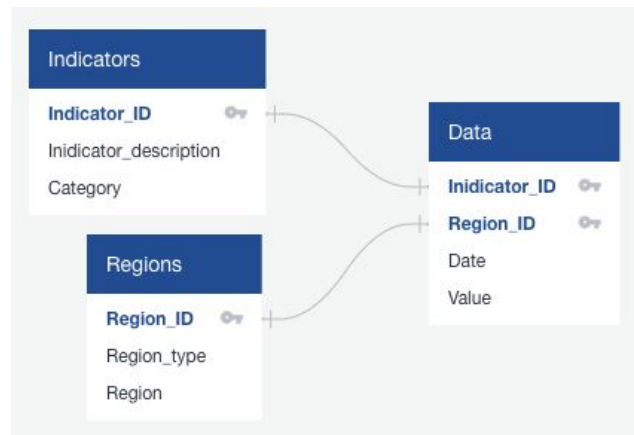
# Challenges

- While the data was geographically based, there was no geographical information provided other than region names and zip codes, so it proved challenging when mapping the data in a visualization
- The region CSV file required some cleaning so that all data lined up, as well as understanding how each region (zip, metro, neigh, city, county) was organized within the "region" column
- Could only call one indicator and one region at a time with API

| | region_id | region_type | region |
|---|---|---|---|
| 0 | 1286 | county | Orange County;CA;Los Angeles-Long Beach-Anahei... |
| 1 | 3175 | county | Philadelphia County;PA;Philadelphia-Camden-Wil... |
| 2 | 3017 | county | Sacramento County;CA;Sacramento-Roseville-Fols... |
| 3 | 401 | county | Bronx County;NY;New York-Newark-Jersey City, N... |
| 4 | 3165 | county | Hillsborough County;FL;Tampa-St. Petersburg-Cl... |

# Transformation(Exploring Data)

- 3 data files extracted from API
  - Indicator
    - Summary of what each indicator_ID represented
      - Property size (Bedroom)
      - Property type (Condo, single family)
    - All ID's tied to category in Home Values and Sales
  - Region
    - Summary of Region ID data based on types
      - City, County, Neighborhood, Metro Area, and Zip-code
  - Data
    - Holds the Home values and dates based on region and indicator ID's

# Transformation(Cleaning)

- Cleaning (PANDAS & PYDANTIC)
  - Compare (MERGE) & drop unnecessary data from Indicator and Region files
  - Utilized PYDANTIC to validate the data types on our extracted files (INT, OBJ)
  - Restructure for analysis: Region file data filtered & split by type(City, County, Metro, Neighborhood, and Zip-Code)
    - Splitting columns & renaming
      - Replacing Null data with 'NA'
      - Editing numerical data (Zip-Codes missing a number - invalid foreign codes, and converting dates)
- Results
  - 28/ 56 Indicators_id's
  - 27,879/ 89,306 Region_id's
  - 1,048,576 Data values w/ Dates

```
1    region_id,region_type,region
2    403211,neigh,"Longwood;NY;New York-Newark-Jersey City, NY-NJ-PA;New York;Bronx County"
```
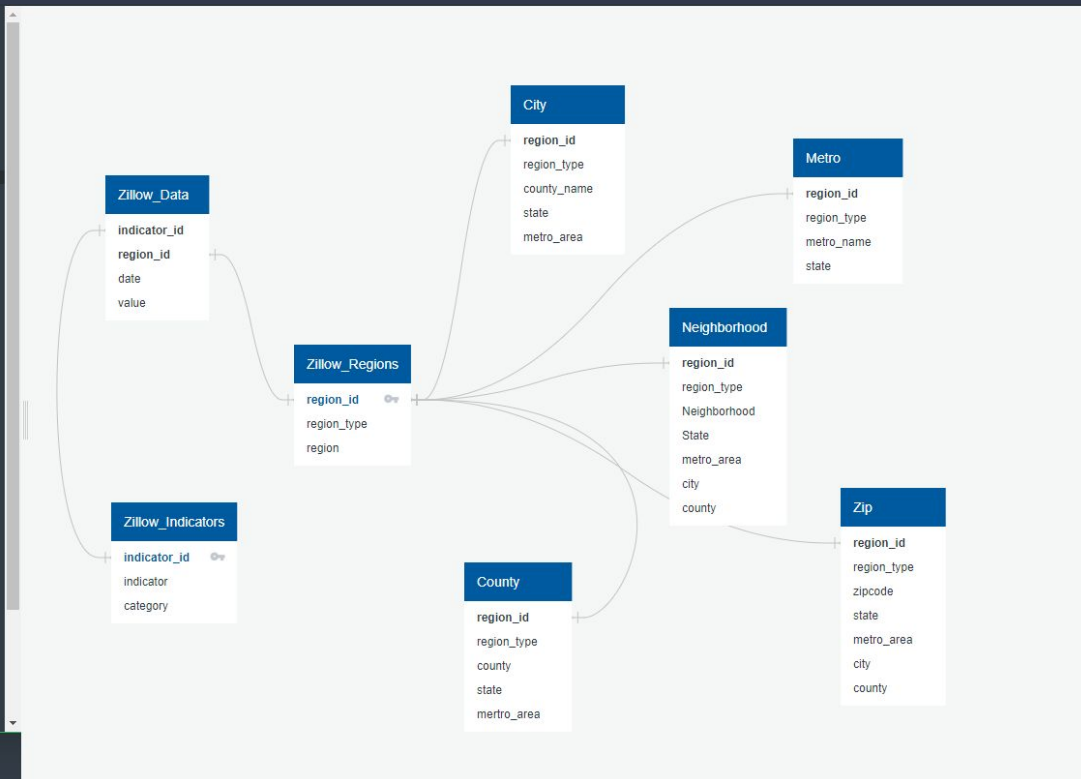
```
region_id,region_type,Neighborhood,State,Metro Area,City,County
403211,neigh,Longwood,NY,"New York-Newark-Jersey City, NY-NJ-PA",New York,Bronx County
```

```
region_id,region_type,region
58924,zip,"2826;RI;Providence-Warwick, RI-MA;Burrillville;Providence County"
```

```
region_id,region_type,Zipcode,State,Metro Area,City,County
58924,zip,02826,RI,"Providence-Warwick, RI-MA",Burrillville,Providence County
```

# ERD Diagram

# SQL Database Storage

```sql
2   CREATE TABLE Indicators(
3       indciator_id VARCHAR PRIMARY KEY NOT NULL,
4       indicator VARCHAR NOT NULL,
5       category VARCHAR NOT NULL
6   );
7
8   --Regions Table
9   CREATE TABLE Regions(
10      region_id INT PRIMARY KEY NOT NULL,
11      region_type VARCHAR NOT NULL,
12      region VARCHAR NOT NULL
13  );
14
15  -- Split Region tables
16  CREATE TABLE City(
17      region_id INT NOT NULL,
18      region_type VARCHAR NOT NULL,
19      City VARCHAR NOT NULL,
20      State VARCHAR NOT NULL,
21      Metro_Area VARCHAR NOT NULL,
22      County VARCHAR NOT NULL,
23      FOREIGN KEY (region_id) REFERENCES Regions(region_id)
24  );
25
26  CREATE TABLE County(
27      region_id INT NOT NULL,
28      region_type VARCHAR NOT NULL,
29      County_Name VARCHAR NOT NULL,
30      State VARCHAR NOT NULL,
31      Metro_Area VARCHAR NOT NULL,
32      FOREIGN KEY (region_id) REFERENCES Regions(region_id)
33  );
34
35  CREATE TABLE Metro(
36      region_id INT NOT NULL,
37      region_type VARCHAR NOT NULL,
38      Metro_name VARCHAR NOT NULL,
39      State VARCHAR NOT NULL,
40      FOREIGN KEY (region_id) REFERENCES Regions(region_id)
```

```sql
40      FOREIGN KEY (region_id) REFERENCES Regions(region_id)
41  );
42
43  CREATE TABLE Neighborhood(
44      region_id INT NOT NULL,
45      region_type VARCHAR NOT NULL,
46      Neighborhood VARCHAR NOT NULL,
47      State VARCHAR NOT NULL,
48      Metro_Area VARCHAR NOT NULL,
49      City VARCHAR NOT NULL,
50      County VARCHAR NOT NULL,
51      FOREIGN KEY (region_id) REFERENCES Regions(region_id)
52  );
53
54
55  CREATE TABLE Zip(
56      region_id INT NOT NULL,
57      region_type VARCHAR NOT NULL,
58      Zipcode INT NOT NULL,
59      State VARCHAR NOT NULL,
60      Metro_Area VARCHAR NOT NULL,
61      City VARCHAR NOT NULL,
62      County VARCHAR,
63      FOREIGN KEY (region_id) REFERENCES Regions(region_id)
64  );
65
66  --Data table
67  CREATE TABLE Data(
68      indicator_id VARCHAR NOT NULL,
69      region_id INT NOT NULL,
70      date DATE NOT NULL,
71      value FLOAT NOT NULL,
72      FOREIGN KEY (indicator_id) REFERENCES Indicators(inidcator_id),
73      FOREIGN KEY (region_id) REFERENCES Regions(region_id)
74  );
75
```

# API Calls

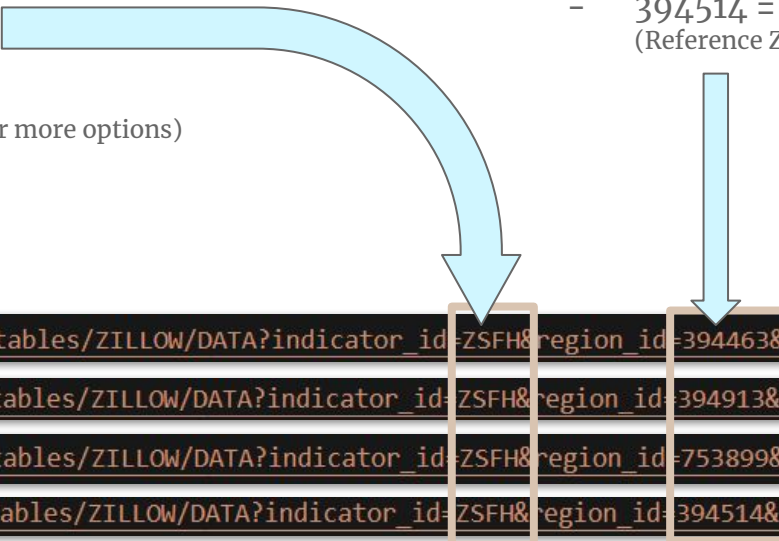**Indicator ID Used:**
- ZSFH: Single Family

**Other Indicator Examples**
- ZATT: All Homes– Top Tier
- ZALL: All Homes
- Z1BR: 1– Bedroom Home
- Z2BR: 2– Bedroom Home

  (Please reference Zillow Indicators.csv for more options)

**Region ID**
- 394463 = Chicago Metro Area
- 394913 = New York City Metro Area
- 753899 = Los Angeles Metro Area
- 394514 = Dallas Metro Area

  (Reference Zillow Regions.csv for more options)

```
url = 'https://data.nasdaq.com/api/v3/datatables/ZILLOW/DATA?indicator_id=ZSFH&region_id=394463&api_key=nhngZKzAkdnohAMb46Kx'
url = 'https://data.nasdaq.com/api/v3/datatables/ZILLOW/DATA?indicator_id=ZSFH&region_id=394913&api_key=nhngZKzAkdnohAMb46Kx'
url = 'https://data.nasdaq.com/api/v3/datatables/ZILLOW/DATA?indicator_id=ZSFH&region_id=753899&api_key=nhngZKzAkdnohAMb46Kx'
url = 'https://data.nasdaq.com/api/v3/datatables/ZILLOW/DATA?indicator_id=ZSFH&region_id=394514&api_key=nhngZKzAkdnohAMb46Kx'
```

# Sample API Call to Create Pandas DataFrame

```python
url = 'https://data.nasdaq.com/api/v3/datatables/ZILLOW/DATA?indicator_id=ZSFH&region_id=394463&api_key=nhngZKzAkdnohAMb46Kx'
response = requests.get(url)
data = response.json()

if response.status_code == 200:
    # Convert response to JSON format
    data = response.json()

    # Extract relevant data
    chicago_data = data['datatable']['data']
    column_names = [column['name'] for column in data['datatable']['columns']]

    # Create DataFrame
    chicago_df = pd.DataFrame(chicago_data, columns=column_names)

chicago_df.head()
```
✓ 0.6s

|   | indicator_id | region_id | date       | value        |
|---|--------------|-----------|------------|--------------|
| 0 | ZSFH         | 394463    | 2024-03-31 | 328937.600843 |
| 1 | ZSFH         | 394463    | 2024-02-29 | 326275.124782 |
| 2 | ZSFH         | 394463    | 2024-01-31 | 324662.193028 |
| 3 | ZSFH         | 394463    | 2023-12-31 | 323830.973942 |
| 4 | ZSFH         | 394463    | 2023-11-30 | 323073.418225 |

```python
chicago_df['date'] = pd.to_datetime(chicago_df['date'])

dates = chicago_df['date']
prices = chicago_df['value']

plt.plot(dates, prices)
plt.title('Chicago House Prices of Single Family Homes')
plt.xlabel('Date')
plt.ylabel('Price ($)')
plt.tight_layout()

# Display the plot
plt.show()
```
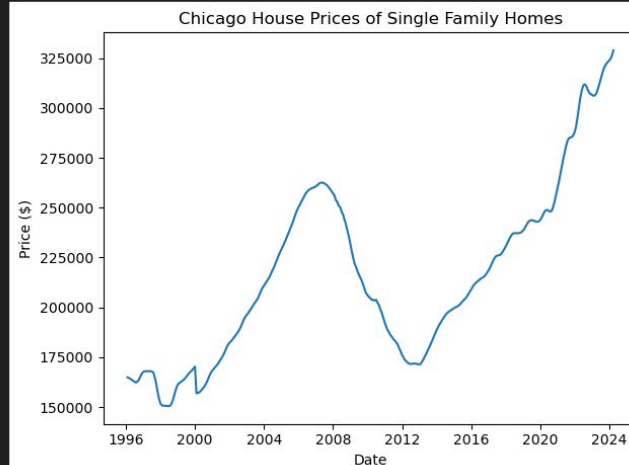✓ 0.0s



Chicago House Prices of Single Family Homes

# API Call with Pandas

```python
# Plot Chicago house prices
plt.plot(chicago_df['date'], chicago_df['value'], label='Chicago')

# Plot New York City house prices
plt.plot(nyc_df['date'], nyc_df['value'], label='New York City')

# Plot Los Angelos house prices
plt.plot(la_df['date'], la_df['value'], label='Los Angelos')

# Plot Dallas house prices
plt.plot(dallas_df['date'], dallas_df['value'], label='Dallas')

# Add titles and labels
plt.title('House Prices of Single Family Homes')
plt.xlabel('Date')
plt.ylabel('Price ($)')
plt.legend()  # Show legend with labels
plt.tight_layout()

# Display the plot
plt.show()
```
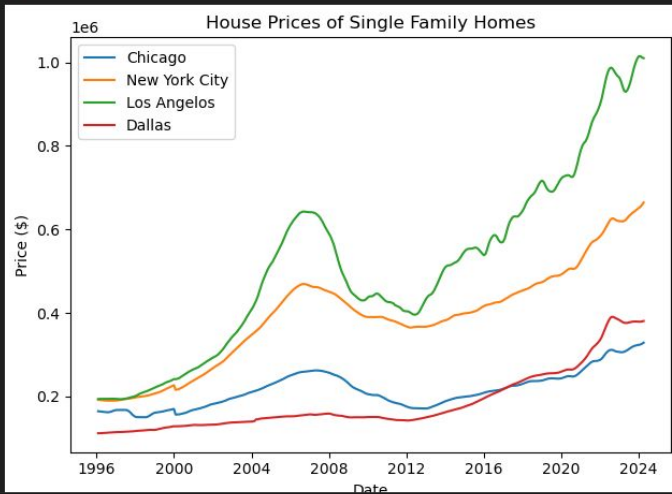✓ 0.3s

By making multiple API calls, each using a distinct region code, we gathered house value data for single-family homes across four major cities. This data was then plotted on a single chart, resulting in four scatterplots illustrating the trend of home values over the years

# GeoPandas

```python
# Read the map of Chicago's community areas
world = gpd.read_file(get_path("geoda.chicago_commpop"))

# Define the bounding box for Chicago [minx, miny, maxx, maxy]
bbox_chicago = [-88, 41.6, -87.3, 42.1]

# Clip the map to the bounding box of Chicago
chicago_clip = world.cx[bbox_chicago[0]:bbox_chicago[2], bbox_chicago[1]:bbox_chicago[3]]

fig, ax = plt.subplots(figsize=(5, 5))
chicago_clip.plot(ax=ax, legend=True)

# Plot the GeoDataFrame on top of the Chicago map
gdf.plot(ax=ax, column='value', markersize=50, cmap='inferno', legend=True)

plt.show()
```
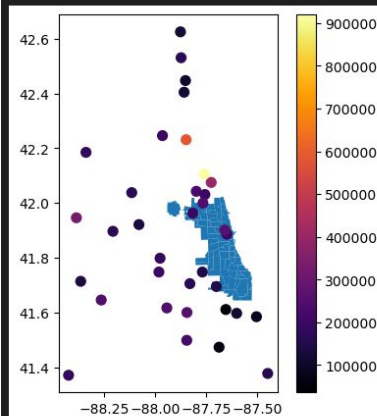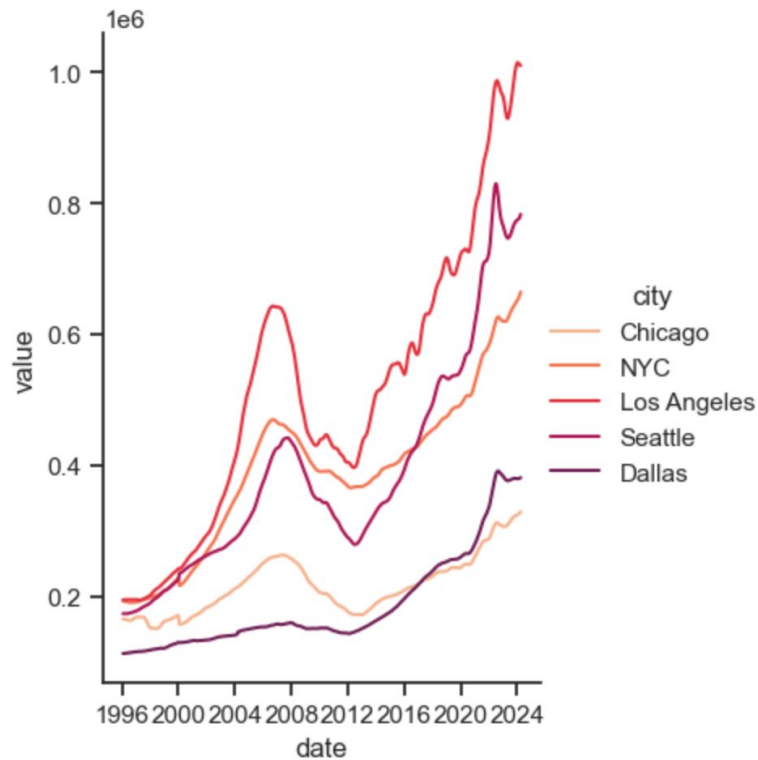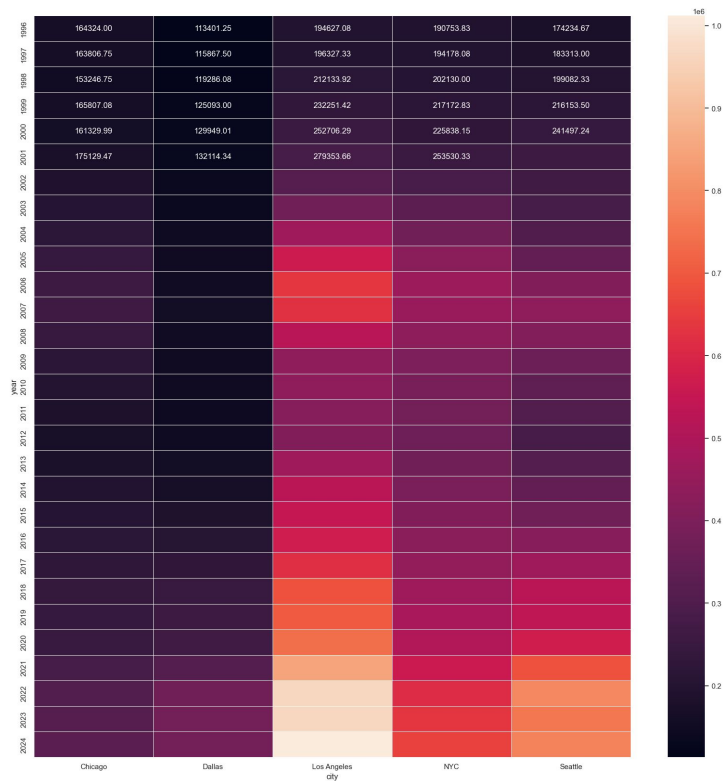✓ 1.1s

```python
# generate latitude and longitude values from zip_codes
latitudes = []
longitudes = []

for zip_code in zip_codes:
    location = geolocator.geocode(zip_code)
    latitudes.append(location.latitude)
    longitudes.append(location.longitude)
```
[150]  ✓  3m 14.7s



|    | region_id | value    | region_type | zipcode | state | state_region                    | city             | county         | latitude  | longitude  |
|----|-----------|----------|-------------|---------|-------|---------------------------------|------------------|----------------|-----------|------------|
| 0  | 78083     | 163152.0 | zip         | 46303   | IN    | Chicago-Naperville-Elgin, IL-IN-WI | Cedar Lake       | Lake County    | 41.378188 | -87.446784 |
| 8  | 78095     | 86862.0  | zip         | 46324   | IN    | Chicago-Naperville-Elgin, IL-IN-WI | Hammond          | Lake County    | 41.585204 | -87.502009 |
| 12 | 81244     | 119609.0 | zip         | 53144   | WI    | Chicago-Naperville-Elgin, IL-IN-WI | Somers           | Kenosha County | 42.624678 | -87.875138 |
| 13 | 81256     | 170217.0 | zip         | 53158   | WI    | Chicago-Naperville-Elgin, IL-IN-WI | Pleasant Prairie | Kenosha County | 42.530150 | -87.871865 |
| 19 | 84308     | 584768.0 | zip         | 60045   | IL    | Chicago-Naperville-Elgin, IL-IN-WI | Lake Forest      | Lake County    | 42.231059 | -87.847344 |

# Seaborn

# Closing Remarks

- Future discovery could be playing with the different indicator IDs to expand search
- Diving into other area regions more in depth
- Researching other market factors to see why certain areas might be priced differently

**Thank you for listening!**

- Questions?