# ETL-Project

**Our Team:**

- Aditya Bhatnagar
- Jonathan Little
- Michael Legg
- Rosario Chiovaro

**Project:**

Comparing real life season statistics of soccer players to their individual ratings in the corresponding FIFA video game year.

**Extract:**

The data was extracted from the source (link below) in the form of CSV's to be later transformed and loaded on a clean database using Pandas and Postgres.

- https://www.kaggle.com/stefanoleone992/fifa-20-complete-player-dataset (CSV)

- https://data.world/cclayford/statbunker-football-statistics (CSV)

**Transform:**

Primarily the data cleaning required the below activities:
   • Selecting the necessary columns
   • Dropping null (NaN) values
   • Changing null (NaN) values to 0

```
In [6]: #transforming dataset 1
        new_players_19_df = players_19_df[['sofifa_id', 'short_name', 'overall', 'pace', 'shooting', 'pass
        ing', 'dribbling', 'defending', 'physic']].copy()
        new_players_19_df = new_players_19_df.dropna()
        new_players_19_df.head()
```

Out[6]:

| | sofifa_id | short_name | overall | pace | shooting | passing | dribbling | defending | physic |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 20801 | Cristiano Ronaldo | 94 | 90.0 | 93.0 | 81.0 | 89.0 | 35.0 | 79.0 |
| 1 | 158023 | L. Messi | 94 | 88.0 | 91.0 | 88.0 | 96.0 | 32.0 | 61.0 |
| 2 | 190871 | Neymar Jr | 92 | 92.0 | 84.0 | 83.0 | 95.0 | 32.0 | 59.0 |
| 4 | 192985 | K. De Bruyne | 91 | 77.0 | 86.0 | 92.0 | 87.0 | 60.0 | 78.0 |
| 5 | 155862 | Sergio Ramos | 91 | 75.0 | 63.0 | 71.0 | 71.0 | 91.0 | 84.0 |

```
In [7]: #transforming dataset 2
        #taking the needed columns
        new_stats_18_19 = stats_18_19[['League', 'Team', 'Season', 'Player', 'Position','Goals','Assists',
        'Yellow Cards']].copy()
        #changing the assist NaN value to 0
        new_stats_18_19.fillna(0,inplace=True)
        #dataset has team totals as a row value, dropping that value
        new_stats_18_19.drop(new_stats_18_19[new_stats_18_19['Player'] == 0].index, inplace = True)


        new_stats_18_19.head()
```

Out[7]:

| | League | Team | Season | Player | Position | Goals | Assists | Yellow Cards |
|---|---|---|---|---|---|---|---|---|
| 1 | Premier League | Manchester City | 2018/19 | Ederson | Goalkeeper | 0 | 0.0 | 2 |
| 2 | Premier League | Manchester City | 2018/19 | Bernardo Silva | Midfielder | 7 | 7.0 | 3 |
| 3 | Premier League | Manchester City | 2018/19 | Aymeric Laporte | Defender | 3 | 0.0 | 3 |
| 4 | Premier League | Manchester City | 2018/19 | Raheem Sterling | Midfielder | 17 | 10.0 | 3 |
| 5 | Premier League | Manchester City | 2018/19 | Sergio Aguero | Forward | 21 | 8.0 | 4 |

```
In [8]: #transforming dataset 3
        new_17_18_df = player_stats_17_18_df[['League', 'Team', 'Season', 'Player', 'Position','Goals', 'A
        ssists', 'Yellow Cards']].copy()
        new_17_18_df.fillna(0, inplace=True)
        new_17_18_df.drop(new_17_18_df[new_17_18_df['Player'] == 0].index, inplace=True)

        new_17_18_df.head()
```

Out[8]:

| | League | Team | Season | Player | Position | Goals | Assists | Yellow Cards |
|---|---|---|---|---|---|---|---|---|
| 1 | Premier League | Manchester City | 2017/18 | Kevin De Bruyne | Midfielder | 8 | 16.0 | 1 |
| 2 | Premier League | Manchester City | 2017/18 | Ederson | Goalkeeper | 0 | 0.0 | 0 |
| 3 | Premier League | Manchester City | 2017/18 | Bernardo Silva | Midfielder | 6 | 0.0 | 0 |
| 4 | Premier League | Manchester City | 2017/18 | Fernandinho | Midfielder | 5 | 0.0 | 7 |
| 5 | Premier League | Manchester City | 2017/18 | Nicolas Otamendi | Defender | 4 | 0.0 | 9 |

```
In [9]: #transforming dataset 4
        new_players_20_df = players_20_df[['sofifa_id', 'short_name', 'overall', 'pace', 'shooting', 'pass
        ing', 'dribbling', 'defending', 'physic']].copy()
        new_players_20_df = new_players_20_df.dropna()
        new_players_20_df.head()
```

Out[9]:

| | sofifa_id | short_name | overall | pace | shooting | passing | dribbling | defending | physic |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 158023 | L. Messi | 94 | 87.0 | 92.0 | 92.0 | 96.0 | 39.0 | 66.0 |
| 1 | 20801 | Cristiano Ronaldo | 93 | 90.0 | 93.0 | 82.0 | 89.0 | 35.0 | 78.0 |
| 2 | 190871 | Neymar Jr | 92 | 91.0 | 85.0 | 87.0 | 95.0 | 32.0 | 58.0 |
| 4 | 183277 | E. Hazard | 91 | 91.0 | 83.0 | 86.0 | 94.0 | 35.0 | 66.0 |
| 5 | 192985 | K. De Bruyne | 91 | 76.0 | 86.0 | 92.0 | 86.0 | 61.0 | 78.0 |

**Load:**

The data was loaded on a SQL database using the procedure shown as below and by creating respective database on PostgresQL

can be used effectively to query the data to compare FIFA video game attributes to real life statistics.

```
In [10]: #loading into the database
         rds_connection_string = "postgres:5Skurlalane!@localhost:5432/Fifa_ETL_Proj"
         engine = create_engine(f'postgresql://{rds_connection_string}')
         new_players_19_df.to_sql(name='fifa_19', con=engine, if_exists='append', index=False)
         new_stats_18_19.to_sql(name='player_stats_18_19',con=engine, if_exists='append', index=False)
         new_17_18_df.to_sql(name='player_stats_17_18', con=engine, if_exists='append', index=False)
         new_players_20_df.to_sql(name='fifa_20', con=engine, if_exists='append', index=False)
```