The background features a series of horizontal color bands: a light tan band at the top left, a dark tan band at the top right, a wide teal band, a light pink band, an orange band, a dark blue band, and a red band at the bottom. Scattered throughout these bands are various organic, hand-drawn shapes in colors matching the bands. The title is centered in the pink band.

IMDB/Bollywood Manipulations and Musings

Chase Khan, Greg Mika

The Raw Datasets

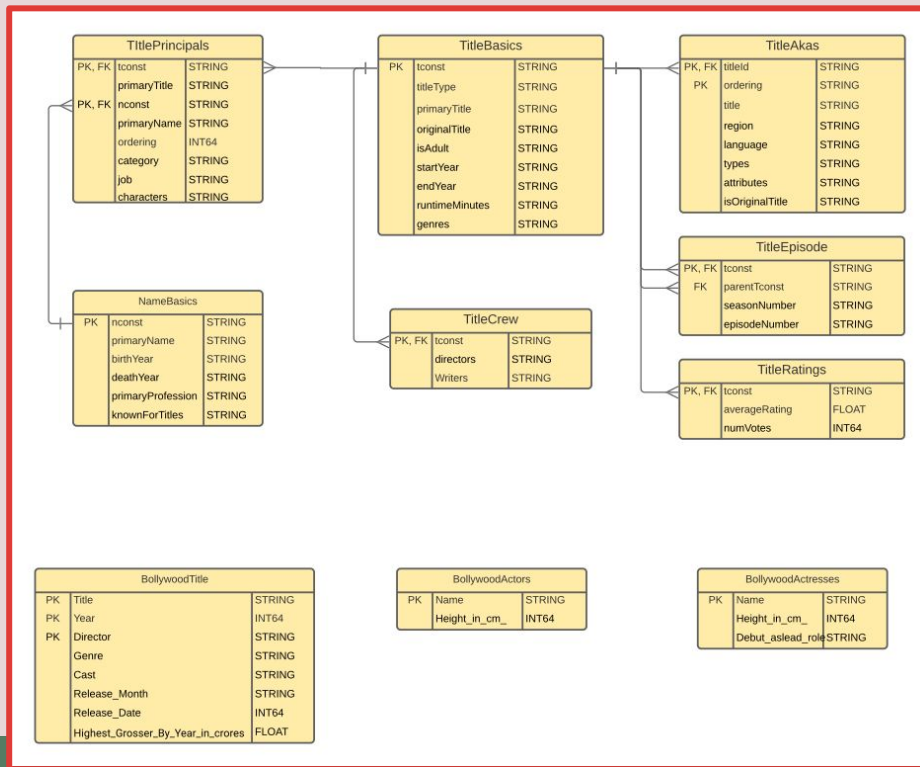
IMDB

- NameBasics
- TitleBasics
- TitleCrew
- TitleEpisode
- TitlePrincipals
- TitleRatings
- TitleAKAs

Bollywood

- Bollywood Actors
- Bollywood Actresses
- Bollywood (Titles)

The Staging Tables



Beam Pipeline Goals

Directors, Writers

Parse a string array

Make a junction table

Join with NameBasics

Bollywood

Create a unique id for the
BollywoodTitles Table

Genres (IMDb and Bollywood), Professions

Parse string arrays

Make a child table

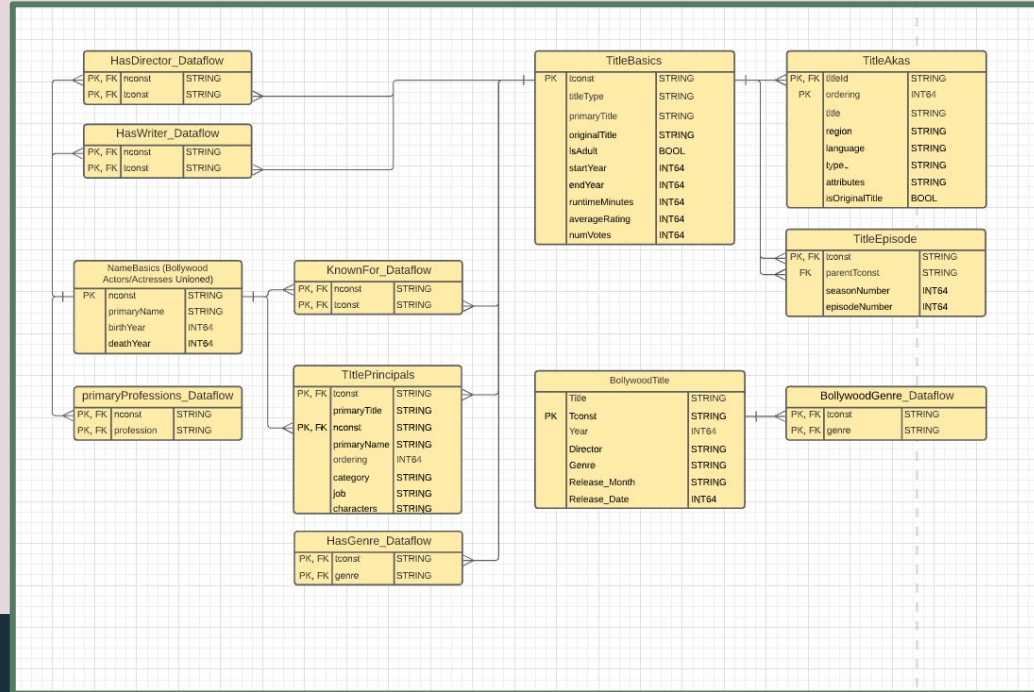
Known For

Parse a string array

Make a junction table

Join with TitleBasics

The Modeled Tables

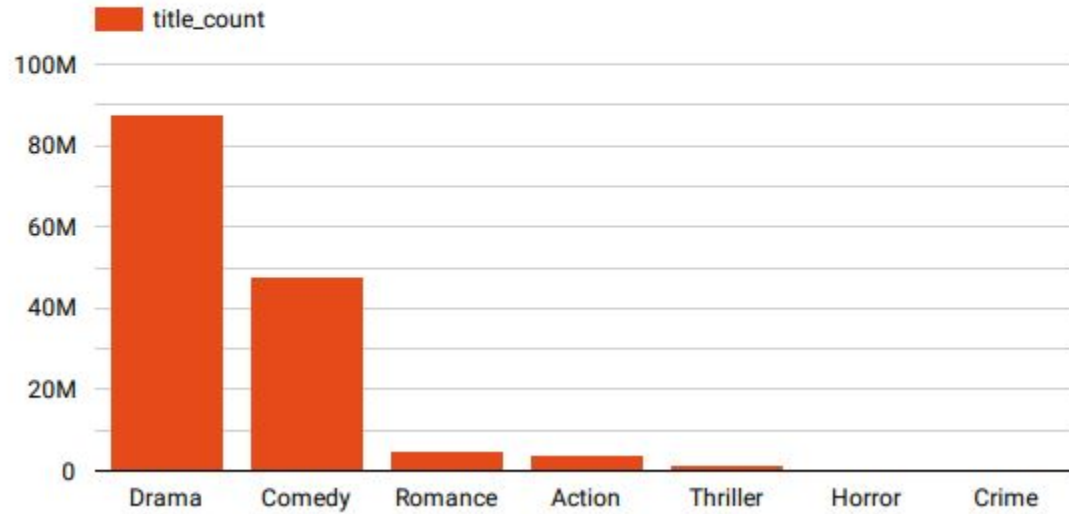




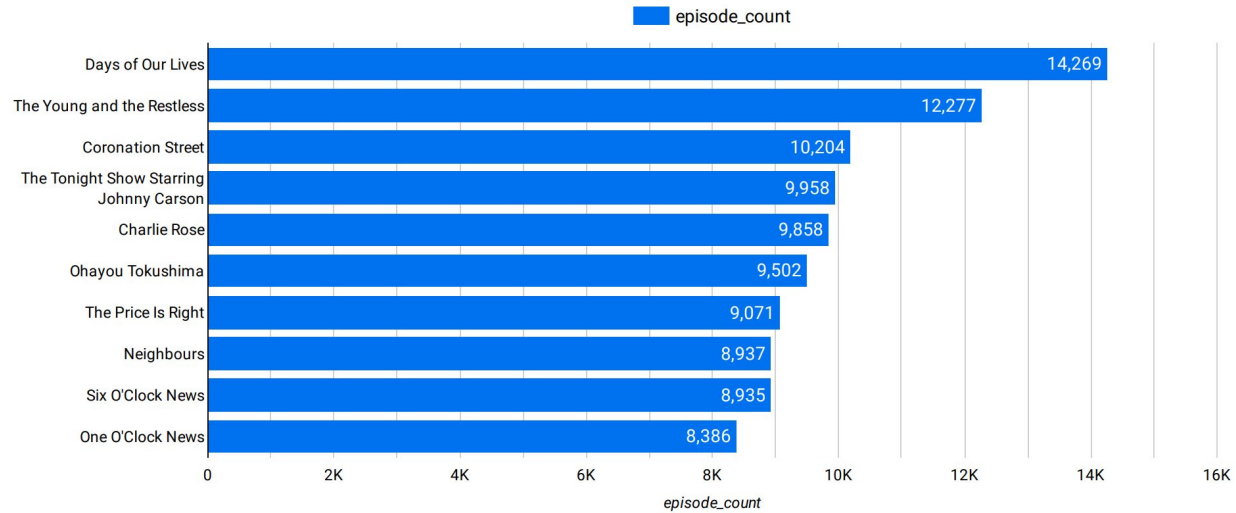
Cross Dataset Queries



Title Count of Genres Represented in Both Datasets



TV Shows with the Most Episodes



Future Improvements

- Write more queries to compare the most popular genre in each dataset, highest rated genres, percentage of titles that make a up genre, etc.
- Integrate the Bollywood titles table into the IMDb title table
 - Find titles that already exist in the IMDb table
 - Add Bollywood titles not in the IMDb table and give them a primary key of the same format as the IMDb table
- Find another secondary data source to fill in gaps in the Bollywood data.

