

# COMP 3380 Final Project

Souvik Ray and Rahul Kumar

## Summary

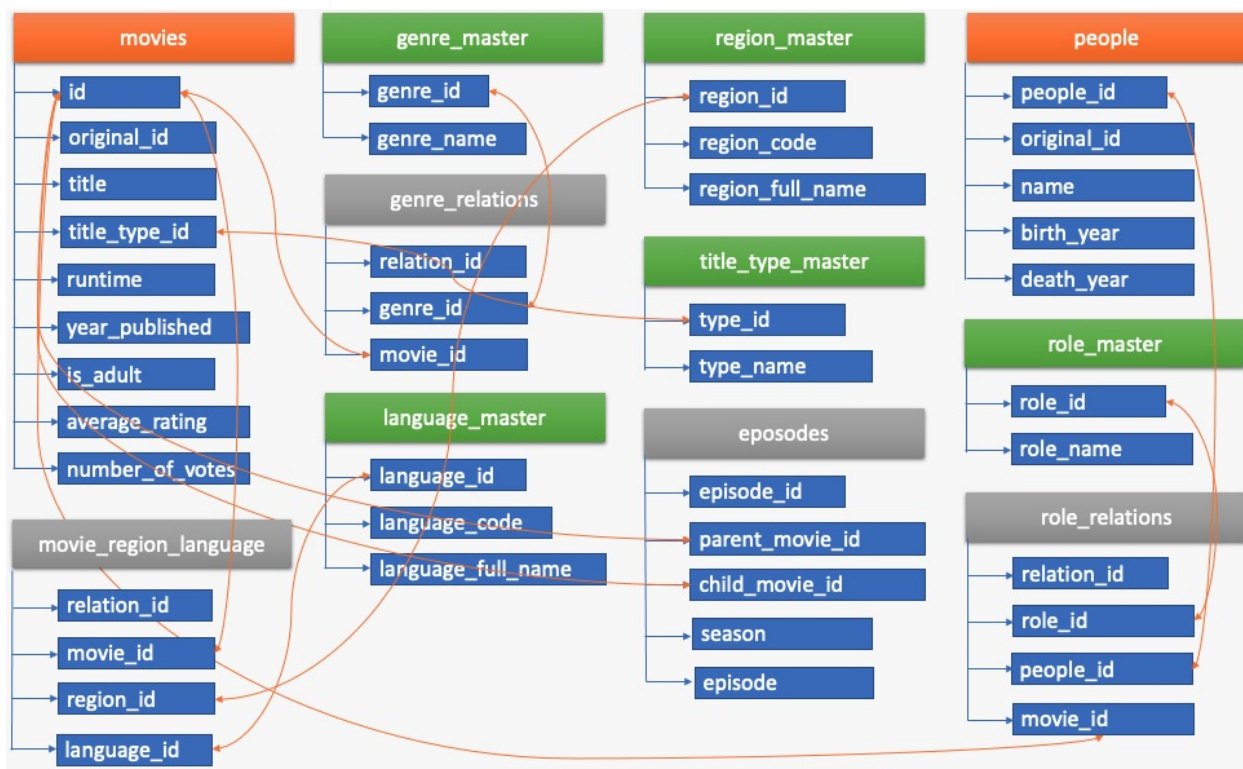
1. This IMDB database was chosen to allow us to access various information which were of different types and could be divided into different tables for quick access of data while having a significant amount of rows to work with.

This database contains 7 different tsv files. The total file size exceeds 5.5 GB and has approximately 12 million data entries in each file.

This link is different than the one previously chosen because the dataset available at the earlier link was removed before we had the opportunity to finish the project.

The link to the dataset - <https://www.imdb.com/interfaces/>

2. The image below gives an overview of how the database is structured.



Movies and People contain most of the information and have been marked in orange. 4 master data tables have been created to give each language, genre, region, title and role a specific id which allows us to reference and parse the data later.

3. This database was created in DB Browser for SQLite. It was populated using the above IMDB datasets and the process was automated using python. The front end UI has also been created using python. To populate the database please use the **upload\_data.py** and uncomment the functions as required at the bottom of the file.

To start the front end UI please use the **data\_analysis.py** file. It might require the user to install several modules like **PyQt5** and **tkinter**. Attached below is an image of how the UI looks like and how it displays results.

List all movie of type (x) can also be exported into a csv file format which automatically creates a new document in the current folder/directory.

Souvik Ray, UoM Project

Choose the Query

List all movies whose rating is

=

1

Find

List all movies of type

short

Export

Find

List all movies where

Adrian Biddle

worked as

actor

Find

List all movies from the region

AD

Find

List all movies with genre

Action

as well as

Action

Find

List number of movies where

someone worked as a Director

between 1900 and 1920

Find

List all Comedy movies where

someone worked between

1890 and 1920

Find

List all movies in which

no one from the movie Conjuring

has ever worked on

Find

List all movies which

has more than 5 rating where

Charles Chaplin worked

Find

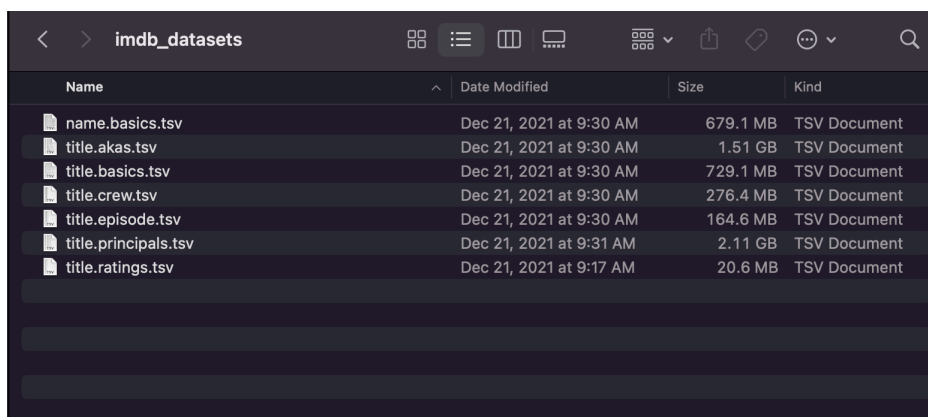
#### 4. Queries —

1. List all movies whose rating is greater than, less than or equal to the selected rating.
2. List all movies of that type.
3. List all movies where person x worked in y role.
4. List all movies from the chosen region.
5. List all movies in which is a part of both genres.
6. List all movies someone worked as a director between certain years.
7. List all comedy movies someone worked between certain years.
8. List all movies where no one from the movie conjuring has ever worked on.
9. List all the movies which has a 5 star rating where Charlie Chaplin has also worked.
10. Count of all the movies of all genres in a given range of years.

#### 5. Links -

Please find below some important links and information -

1. The entire project folder link - <https://1drv.ms/u/s!AhVTjYyRL-ywpzJS51po3itZACxW?e=RnIBG3>
2. Dataset - <https://www.imdb.com/interfaces/>
3. The datasets must have the same name as they were when they are downloaded from IMDB website. They must be kept in a single folder called 'imdb\_datasets'.
4. Due to file size issues we



Name	Date Modified	Size	Kind
name.basics.tsv	Dec 21, 2021 at 9:30 AM	679.1 MB	TSV Document
title.akas.tsv	Dec 21, 2021 at 9:30 AM	1.51 GB	TSV Document
title.basics.tsv	Dec 21, 2021 at 9:30 AM	729.1 MB	TSV Document
title.crew.tsv	Dec 21, 2021 at 9:30 AM	276.4 MB	TSV Document
title.episode.tsv	Dec 21, 2021 at 9:30 AM	164.6 MB	TSV Document
title.principals.tsv	Dec 21, 2021 at 9:31 AM	2.11 GB	TSV Document
title.ratings.tsv	Dec 21, 2021 at 9:17 AM	20.6 MB	TSV Document

decided to not include the datasets in the submitted folder. For every program to work - please download and keep all datasets in the supplied folder named - 'imdb\_datasets'.