



Blavatnik School of Computer Science  
Tel Aviv University

# Database Systems Course

## 0368.3458

### Home Work 3 – Movies Web App

#### Software Docs

By:  
Sidney Feiner  
Omer Moses

## Contents

Software Documentation .....	3
DB scheme structure .....	3
DB optimization performed.....	4
Queries Description .....	5
Get Movies.....	5
Get Crew .....	5
Get Cast .....	6
Multi Role Actors .....	7
Lookalike Movies .....	7
Best Profit Members.....	8
Loyal Crew Members.....	9
Code Structure.....	10
Description Of The API.....	10
Used Endpoints.....	10
Provided Endpoint .....	10
/cast .....	10
/crew.....	11
/misc .....	11
/movies .....	13
Opening The Swagger guidelines .....	14
General flow of the application.....	14

## Software Documentation

### DB scheme structure

The design of the schema will be presented in several stages.

The key building blocks of our information are:

- movies
- crew
- cast

These building blocks contain the most basic information that needs to be used to generate more complex queries that incorporate information that users will want to know.

movies	
id	int
title_id	int
original_title_id	int
original_language	varchar(2)
popularity	decimal(9,3)
release_date	date
tagline	varchar(500)
is_adult	tinyint(1)
overview	varchar(1500)
budget_usd	int
revenue_usd	bigint
runtime_minutes	int
status_id	int
vote_avg	decimal(4,2)
vote_cnt	int

crew	
id	int
person_id	int
gender_id	int
job_id	int
department_id	int
movie_id	int
id_in_crew	int

cast	
id	int
person_id	int
character_name_id	int
order	int
cast_id	int
movie_id	int
id_in_cast	int

For information that does not change frequently, or that is not unique to a movie, actor, or team member, we have created dedicated tables. For example, if a staff member changes his department, and interested in getting offers in that new department, we can easily change that and it will be updated automatically. This way we are able to maintain scalable and efficient data updating.

countries	
id	int
country_name	varchar(50)
country_name_iso_3166_1	varchar(2)

languages	
id	int
language	varchar(20)
language_iso_639_1	varchar(2)

production_companies	
id	int
production_company	varchar(200)

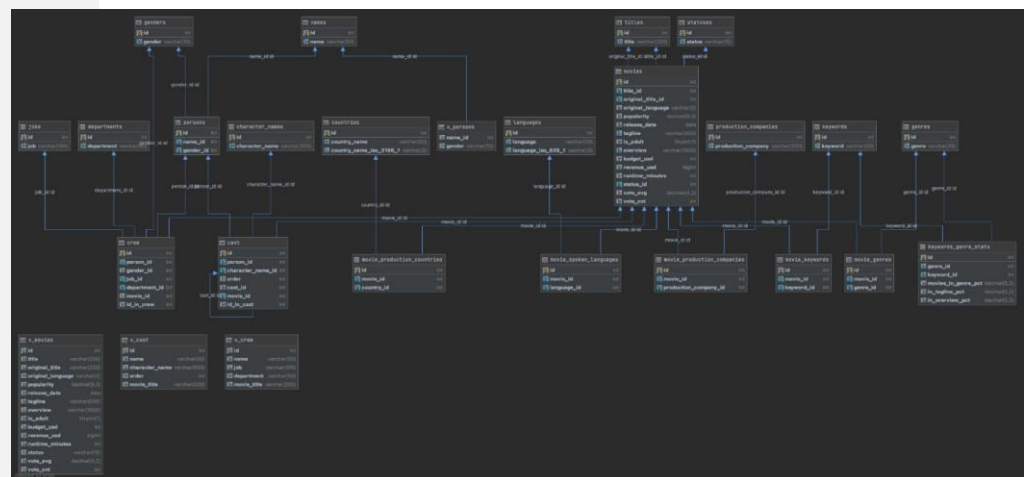
genres	
id	int
genre	varchar(20)

jobs	
id	int
job	varchar(100)

departments	
id	int
department	varchar(100)

Finally, connections between records are made based on key and foreign keys. So, for instance, if we take movie XXX his spoken language is stored in “movie\_spoken\_languages” under some foreign key – language\_id which is a key at “languages” table. Clearly, this information is not readable by the user. To solve this problem, we have created views that present the data in a textual way. Thus, for the case we presented, will be displayed to the user when executing a query, the English language for the movie XXX.

Our full DB schema:



### DB optimization performed

In order to provide a positive user experience, along with optimal utilization of our resources, we have performed a number of DB optimizations:

1. Cron based data update – In order to keep our data fresh, we run an update across all of our data kept in our DB. Since the data isn’t updated frequently, we will run a procedure that updates the data once a day. Note that during these procedures the user will have the ability to use the app based on old data, the data is updated atomically.
2. Indices – Our DB schema is build based on the features we provide to our users. After defining the features, we designed the schema to support those feature using indexes on data we would like to get easily. The following indices are applied:

Table Name	Index Over
names	name
keywords	keyword
character_names	character name

<b>jobs</b>	jobs
<b>departments</b>	department
<b>genres</b>	genre
<b>titles</b>	title
<b>production_companies</b>	Production companie
<b>countries</b>	country name country name iso
<b>languages</b>	language language iso
<b>statuses</b>	status

### Queries Description

The following section describes our main implemented queries.

#### Get Movies

The query is an abstraction of a basic query pulling data from “movies” table.

Formally we can look at this query as

```
SELECT {projection choice}
FROM movies_table
WHERE {user choice}
```

#### Notes:

- Leaving the parameter projection choice empty is equal to SELECT \*
- Leaving the user choice empty returns data from the table without conditioning

#### Optimization:

- Indices

#### How does the DB support this query:

- As mentioned earlier, movies table has a view the represents user readable data.
- Movies are the basis for data, and are used by us for the more complex queries built on the basis of these movies, the information is arranged in a separate basic table. In order to create an identification on the films they get a unique ID on the basis of which we can build more complex queries

#### Get Crew

The query is an abstraction of a basic query pulling data from “crew” table.

Formally we can look at this query as

```
SELECT {projection choice}
```

```
FROM crew_table
WHERE {user choice}
```

Notes:

- Leaving the parameter projection choice empty is equal to SELECT \*
- Leaving the user choice empty returns data from the table without conditioning

Optimization:

- Indices

How does the DB support this query:

- As mentioned earlier, movies table has a view the represents user readable data.
- Crew is one of the basis for data, and used by us for the more complex queries built on this basis of these crew members, the information is arranged in a separate basic table. In order to create an identification on the members, they get a unique ID on the basis of which we can build more complex queries.

Get Cast

The query is an abstraction of a basic query pulling data from “cast” table.

Formally we can look at this query as

```
SELECT {projection choice}
FROM crew_table
WHERE {user choice}
```

Notes:

- Leaving the parameter projection choice empty is equal to SELECT \*
- Leaving the user choice empty returns data from the table without conditioning

Optimization:

- Indices

How does the DB support this query:

- As mentioned earlier, movies table has a view the represents user readable data.
- Cast is one of the basis for data, and used by us for the more complex queries built on this basis of these crew members, the information is arranged in a separate basic table. In order to create an identification on

the members, they get a unique ID on the basis of which we can build more complex queries.

#### Multi Role Actors

In search for talented actors the query return actors that played multiple roles in a single movie.

Gender can be given and filtered by, but is optional.

```
1 select n.name, t.title, group_concat(distinct cn.character_name) as characters
2 from cast
3 join names n on cast.name_id = n.id
4 join movies m on cast.movie_id = m.id
5 join character_names cn on cast.character_name_id = cn.id
6 join titles t on m.original_title_id = t.id
7 (where_clause)
8 group by 1, 2
9 having count(distinct cn.character_name) >= 2
```

Optimization:

- Indices

How does the DB support this query: this query uses Cast and Movies tables and searches based on name\_id which is a key in Cast table and movies\_id and title which is foreign key from Movies table.

#### Lookalike Movies

The query takes a movie ID and returns lookalike movies. We defined lookalike movies are been graded based on

- Genre
- mutual crew/cast members

Eventually the query sorts the results based on the movie's reassembles score.

```

1 with same_genres as
2 (
3     select mg1.movie_id as movie_id_1,
4         mg2.movie_id as movie_id_2,
5         count(case when mg1.genre_id = mg2.genre_id then 1 end) as cnt
6     from (contract.MOVIE_GENRES_TABLE) mg1
7     join movie_genres mg2 on mg1.movie_id <> mg2.movie_id
8     where mg1.movie_id = %(movie_id)s
9     group by 1, 2
10 ),
11 same_cast as (
12     select c1.movie_id as movie_id_1,
13         c2.movie_id as movie_id_2,
14         count(case when c1.name_id = c2.name_id then 1 end) as cnt
15     from (contract.CAST_TABLE) c1
16     join (contract.CAST_TABLE) c2 on c1.movie_id <> c2.movie_id
17     where c1.movie_id = %(movie_id)s
18     group by 1, 2
19 ),
20 same_crew as (
21     select c1.movie_id as movie_id_1,
22         c2.movie_id as movie_id_2,
23         count(case when c1.name_id = c2.name_id then 1 end) as cnt
24     from (contract.CREW_TABLE) c1
25     join (contract.CREW_TABLE) c2 on c1.movie_id <> c2.movie_id
26     where c1.movie_id = %(movie_id)s
27     group by 1, 2
28 )
29 select m1.id as requested_movie_id,
30     t1.title as requested_movie,
31     m2.id as movie_id,
32     t2.title as title,
33     m2.original_language,
34     sg.cnt as mutual_genres_amt,
35     sc.cnt as mutual_cast_amt,
36     scr.cnt as mutual_crew_amt,
37     m2.vote_avg,
38     IF(m1.original_language = m2.original_language, 5, 0) + 3 * sg.cnt + 2 * sc.cnt +
39     scr.cnt as lookalike_score
40 from (contract.MOVIES_TABLE) m1
41 join movies m2 on m1.id <> m2.id
42 join same_genres sg on sg.movie_id_1 = m1.id and sg.movie_id_2 = m2.id
43 join same_cast sc on sc.movie_id_1 = m1.id and sc.movie_id_2 = m2.id
44 join same_crew scr on scr.movie_id_1 = m1.id and scr.movie_id_2 = m2.id
45 left join titles t1 on m1.title_id = t1.id
46 left join titles t2 on m2.title_id = t2.id
47 where m1.id = %(movie_id)s
48 order by lookalike_score desc,
49     m2.vote_avg desc

```

Optimization:

- Indices

How does the DB support this query: the query searches based on genre that is indexed in his table for quicker search. The query also uses mutual cast members and crew workers that are primary keys in their table

### Best Profit Members

The query takes a movie net revenue and divides it by the number of the entire workers who took place creating this movie (crew + actors).

```

1 with movie_members_amt as (
2     select movie_id, count(distinct name_id) as members
3     from (
4         select movie_id, name_id
5         from cast
6         union all
7         select movie_id, name_id
8         from crew
9     ) a
10    group by 1
11 )
12 select m.id as movie_id, t.title as movie_name, (m.revenue_usd - m.budget_usd) / mm.members as profit_rate
13 from movies m
14 join movie_members_amt mm on m.id = mm.movie_id
15 join movie_genres mg on m.id = mg.movie_id
16 join titles t on m.title_id = t.id
17 where genre_id = 18 and m.revenue_usd is not null and m.budget_usd is not null
18 order by 1, profit_rate desc

```

Optimization:

[SF1] עם הערות: Index on movie\_id in all relevant tables



- Indices

How does the DB support this query: Each movie registry in Movies table, has the net revenue value and this query gets the movie id as input. Next, movie\_id is a foreign key in both Crew and Cast table so we can use this index to search for relevant registries.

#### Loyal Crew Members

The query searches for crew members who worked for single production company – showing their loyalty to their company.

```

1 with crew_company as (
2   select crew.name_id, crew.job_id, mpc.production_company_id
3   from crew
4   join movie_production_companies mpc on crew.movie_id = mpc.movie_id
5 )
6 select n.name, pc.production_company, group_concat(distinct jobs.job) as jobs
7 from crew_company cc1
8   join names n on cc1.name_id = n.id
9   join production_companies pc on cc1.production_company_id = pc.id
10  join jobs on cc1.job_id = jobs.id
11  where exists(select 1 from crew where cc1.name_id = crew.name_id (jobs_filter))
12  and not exists(select 1
13                  from crew_company cc2
14                  where cc1.name_id = cc2.name_id
15                    and cc1.production_company_id <> cc2.production_company_id)
16 group by 1, 2

```

Optimization:

- Indices

How does the DB support this query: Each crew member registry in Crew table, holds a foreign key from movies\_productions\_companies called movie\_id. Using this movie id we can easily find the correct movies registry and take the production company id which is a foreign key in movies\_productions\_companies table.

## Code Structure

### CREATE-DB-SCRIPT

Contains the SQL queries we used to create the tables, tables and indices in our database.

### API-DATA-RETRIEVE

This file is used to get the data from our API and csv files, parse it and propagate to our tables.

First, the system gathers the data and parse it from both of our sources.

Next, the data is being wrapped in a generator getting ready to be loaded into the database.

### SERVER

This package contains services that responsible to:

- Querying the data at aql.py
- Exposing the API the our users will use while using the app at server.py

## Description Of The API

### Used Endpoints

We used themoviedb API to enrich our data base. Using rest API get call we get from the endpoint json formatted response with our DB's data.

First, we created a generator that iterates of movies data held in a dictionary. Using this generator we populate the data into the tables that we showed in our schema above.

### Provided Endpoint

Our app users may choose to use any of our 9 provided endpoints. For the user's convenience, we provide a swagger that explains each endpoint, such as the type of operation, parameters, sample query and exceptions. You may use this [link](#) for opening the swagger guidelines.

/cast

### Get Cast

The user may choose any of the following parameters he would like to know about the cast stored in our DB.

Note that leaving the parameters empty will return the all columns.

Parameter:

Name	Type	Description
<b>ID</b>	integer	ID of the member
<b>name</b>	string	Name of the member
<b>character_name</b>	string	character name of cast member in this movie

<b>gender</b>	string	gender
<b>order</b>	integer	order in the credits of cast member for this specific movie
<b>movie_title</b>	string	title of movie this member was cast in
<b>limit</b>	integer	If given, will limit the amount of returned records

#### Get Cast Multirole

Return actors that played multiple roles in a single movie. The user may choose one of the following parameters to filter the results getting from the app.

Parameter:

Name	Type	Description
<b>Gender</b>	string	gender
<b>Limit</b>	integer	If given, will limit the amount of returned records

/crew

#### Get Crew

The user may choose any of the following parameters he would like to know about the crew stored in our DB.

Note that leaving the parameters empty will return the all columns.

Parameter:

Name	Type	Description
<b>ID</b>	integer	ID of the member
<b>Name</b>	string	Name of the member
<b>gender</b>	string	gender
<b>job</b>	string	The Role of the member
<b>department</b>	string	Department of the member
<b>movie_title</b>	string	title of movie this member was cast in
<b>limit</b>	integer	If given, will limit the amount of returned records

/misc

#### Get Best Profit Per Worker

For a given genre, this endpoint returns the movie with the best profit per worker, where best profit per worker is the result of dividing the net revenue of the movie by all workers (cast + crew). Eventually, the result is movies where every worker had in the

average, the biggest contribution to the movie's profit.

Parameter:

Name	Type	Description
genre	string	ID of the member
limit	integer	If given, will limit the amount of returned records

#### Get Genre Distribution

Returns the distribution of genres for a given query. If no query is given, all movies will be taken into account.

Parameter:

Name	Type	Description
genre	enum	Use BOOLEAN QUERY syntax, which is described here: <a href="https://dev.mysql.com/doc/refman/8.0/en/fulltext-boolean.html">https://dev.mysql.com/doc/refman/8.0/en/fulltext-boolean.html</a>

#### Get Keyword Genre Stats

Given genres, returns the most common keywords for these genres, with extra data regarding in how many movies of this genre (in percentage), the keyword appears in its description (either tagline or overview). More than limit keywords can be returned per genre, if some of these keywords have the same weight in the genre.

Note that leaving the parameters empty will return the all columns.

Parameter:

Name	Type	Description
ID	integer	ID of the member
name	string	Name of the member
gender	string	gender
job	string	The Role of the member
department	string	Department of the member
Movie Title	string	title of movie this member was cast in
limit	integer	If given, will limit the amount of returned records

#### Get Loyal Crew Members

The user may find crew members that work only with a single production company.

API returns data for a specific job (or jobs). If no job is given, all jobs will be searched.

Parameter:

Name	Type	Description
------	------	-------------

<b>job</b>	array of comma separated strings	List of jobs
<b>limit</b>	integer	If given, will limit the amount of returned records

#### [/movies](#) Get Movies

The user may choose any of the following parameters he would like to know about the cast stored in our DB.

Note that leaving the parameters empty will return the all columns.

Parameter:

Name	Type	Description
<b>id</b>	integer	ID of the movie
<b>title</b>	string	Name of the movie
<b>original_language</b>	string	
<b>popularity</b>	float	
<b>release_date</b>	string	Formatted as yyyy-mm-dd
<b>tagline</b>	string	
<b>overviews</b>	string	
<b>budget_usd</b>	integer	
<b>revenue_usd</b>	integer	
<b>runtime_minutes</b>	integer	
<b>status</b>	string	
<b>vote_avg</b>	float	On a scale of 0-10
<b>vote_cnt</b>	int	
<b>projection</b>	string	array of columns that we want the API to return. Can be an array containing any of the above columns. If no projection is given, all columns will be returned
<b>limit</b>	integer	If given, will limit the amount of returned records

#### Get Movies Lookalike

This query takes a movie ID and returns lookalike movies.

We defined lookalike movies are been graded based on

- Genre
- mutual crew/cast members

Results are sorted by the movie's reassembles score.

parameters:

Name	Type	Description
<b>id</b>	integer	ID of the member
<b>name</b>	string	Name of the member
<b>gender</b>	string	gender
<b>job</b>	string	The Role of the member
<b>department</b>	string	Department of the member
<b>movie_title</b>	string	title of movie this member was cast in
<b>limit</b>	integer	If given, will limit the amount of returned records

### Opening The Swagger guidelines

Open the file SERVER/run-server.sh and insert correct user/password.

Once done, run this script and then the server will start running (on port 5050) (and you'll be able to use all our APIs) and in addition, you can see the whole API documentation at the endpoint: `SERVER:5050/apidocs/`.

For example, if you're running locally: `http://localhost:5050/apidocs/`

### General flow of the application

This section will provide three main flows that users experience while using the app.

Both on the user's side and the app side we should know the role of the user in the cinema world, therefore the first flow we will registration. Next goes into our login flow where he should have successful identification, passing that, the landing page will be depended on the user's role. Eventually the user may move across our app features getting the data he need.

#### Registration Flow

The following form should be filled by the user. Note that this is a free account. As future plans we would like to create subscribers account:



#### Login Flow

The following basic flow chart describe the login flow stats machine.

