# Design

*"Design a distributed database based on your analysis part's insights."*

## Introduction

The main principles which guided us in the design process of the distributed data base were as followed:

1. We assumed the database will contain more read operations than write operations from the users, therefore, we'll design the database with no restriction on having duplicate data in different fragments and sites.
2. The fragmentation will try to place tuples as close to the users (geographically) which query the tuples the most. This is import since we could have duplicated all the data. We want to consider the trade-off between duplicating the data to keeping the database as simple as possible.
   (Note: Users from 'Kibuts Gesher' are closest to 'Tiberias' site.)
3. The users output from the queries is the movie ID. All the additional description as: title, overview etc… never queried about, is irrelevant, allowing us to filter columns out of the database.
4. The main guidance we were given as mentioned by the staff:
   " Our obligation to you is – the only queries to be asked are the queries as they are mentioned in the data, your obligation to us is to be able to answer theses queries and return the movies ID. All redundant data can be deleted. "

## Part 1 – Join

The two tables of data named: movies.csv and credits.csv, are all the data we need to take into consideration in the design section. For the ease of the design process, we decided to refer to them as one relation named R.

To achieve this, we can Join both tables on the 'movie_id' column.

Header of R:

| movie_id | genres | overview | Production _companies | Production _countries | Release _date | revenue | Spoken _languages | tagline | title | cities | cast | crew |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

## Part 2 – Filter by Column (Vertical)

First, we'll filter all the columns never required by the users nondependent on our insights.

'overview', 'revenue', 'tagline' and 'title' are the columns we delete from the database.

Reminder: based on the user queries we cleaned the data of many columns, for example: from 'crew' and 'cast' columns, we left only the actors and director of the movies.

After the filter process we have the following header for R: (renamed 'crew' and 'cast')

| movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | cities | actors | director |
|---|---|---|---|---|---|---|---|---|

Secondly, we'll filter all the columns never required by the users as found in our insights per city.

The relevant insights are 4.1 and 4.2 from the previous question.

From insight 4.1, we can remove the column 'actors' for the data stored on the Eilat site. Users from Eilat never query regarding actors:

Resulting in the header $R_{Eilat}$:

| movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | cities | director |
|---|---|---|---|---|---|---|---|

From insight 4.2, we can use the users behaviour regarding querying directors, as a feature for the basis of a fragmentation.

As opposed to the fragmentation we concluded from insight 4.1 for Eilat's site, in this case we cannot remove the director column permanently from the site which query directors only up to 4% of their queries.('Tiberias', 'Tel Aviv', 'Jerusalem', 'Eilat') The reason being, if we remove the column, believing we keep the director column only in another site ('Haifa'), once we continue to filter the data from the site containing the director column - we may remove useful tuples as a result of the sites constraints, without taking all the other sites into consideration. (Example: in Haifa site, we will filter movies from 2010 and below, while others site query from 1990 onwards. All the tuples in between with their directors would be lost)

Therefore, we came up with the following solution:

We'll create a separate table, containing only: 'movie_id' as a key and 'director'. This table wont be duplicated on each site, because of small use as we showed in insight 4.2. This table is supposed to answer the queries of 4 sites which have low demand for director information, will be stored on one of them and will share the data when requested by the other sites.

We chose to store this table on Jerusalem's site, because the users from Jerusalem query the data most (Jerusalem – 30,873 queries, 'Tel Aviv' – 24,839, 'Eilat' – 17,348, 'Tiberias' - 5288). All 4 sites query a director ~4% of the time, therefore, Jerusalem's number of queries is greater than others.

The director table will contain all the movies from 1990 onwards (derived from insight 2, explained further in the following pages.)

We remove the column 'director' for the data stored on all sites except Haifa. Resulting in the headers:

$R_{Tiberias}$:

| movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | Cities | actors |
|---|---|---|---|---|---|---|---|

$R_{Tel-Aviv}$:

| movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|

$R_{Jerusalem}$:

| movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|

| movie_id | director |
|---|---|

$R_{Eilat}$:

| movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | cities |
|---|---|---|---|---|---|---|

And $R_{Haifa}$:

| movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | cities | actors | director |
|---|---|---|---|---|---|---|---|---|

## Part 3 – Filter by Row (Horizontal)

From insight 1 in the previous question, we can divide horizontally each sites data, based on the location in which the movies take place. If a site's users never query regarding to certain cities, the relevant movies for the site are those they are going to query about. See figure 1.2 in the insights for the full thought process behind our conclusion.

The result from insight 1 on the fragmentation:

$R_{Tiberias}$:

| Constraints from… | movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Tiberias, Haifa | |

$R_{Tel-Aviv}$:

| Constraints from… | movie_id | genres | Production_companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Jerusalem, Tel-Aviv | |

$R_{Jerusalem}$:

| Constraints | movie_id | director |
|---|---|---|
| Movie release date 1990 onwards | | |

| Constraints from… | movie_id | genres | Production _companies | Production_countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Jerusalem, Tel-Aviv | |

$R_{Eilat}$:

| Constraints from… | movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | cities |
|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Eilat |

$R_{Haifa}$:

| Constraints from… | movie_id | genres | Production _companies | Production _countries | Release _date | Spoken _languages | cities | actors | director |
|---|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Tel-Aviv, Haifa, Tiberias | | |

From insight 2 in the previous question, we can divide horizontally each sites data, based on the release date of the movies each sites users query. We found that a vast number of movies are never queried about as a result of users requesting a minimum release date for the movies. See figure 2.1 in the insights section.

The result from insight 2 on the fragmentation:

$R_{Tiberias}$:

| Constraints from… | movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Tiberias, Haifa | |
| Insight 2 | | | | | 1990 - current | | | |

$R_{Tel-Aviv}$:

| Constraints from… | movie_id | genres | Production _companies | Production_countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Jerusalem, Tel-Aviv | |
| Insight 2 | | | | | 2010 - current | | | |

$R_{Jerusalem}$:

| Constraints | movie_id | director |
|---|---|---|
| Movie release date 1990 onwards | | |

| Constraints from… | movie_id | genres | Production _companies | Production_countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Jerusalem, Tel-Aviv | |
| Insight 2 | | | | | 1990 - current | | | |

$R_{Eilat}$:

| Constraints from… | movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | cities |
|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Eilat |
| Insight 2 | | | | | 1990 - current | | |

$R_{Haifa}$:

| Constraints from… | movie_id | genres | Production _companies | Production _countries | Release _date | Spoken _languages | cities | actors | director |
|---|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Tel-Aviv, Haifa, Tiberias | | |
| Insight 2 | | | | | 2010 - current | | | | Insight 2 |

We debated in length the way we would be able to use insight 3 for the dividing of the data based on the genres. At first, we tried dividing the sites where there was a preferred genre in the queries by location (see appendix 1 for the design reached).

The problem with this was that the genre column in the query data was multi variable, meaning – if 'Drama' is the most popular, but it always appears with other genres, then dividing the data based on the appearance of the genre wouldn't contribute to trying to keep fragmentations oriented towards the queries. We would need to anyway need to iterate over the fragment and the total data every query.

Second, we tried finding some connection between the popular genres found from the queries data by location and the genres in the movies data, which is also a multi variable column. We were happy to find that 37.4% of the movies are described with only 1 genre. We further checked the most popular genre within these movies, finding 'Drama' appearing 5,000 times as a single genre.

Third, we investigated the number of queries made by users from Haifa, which only query for 'Drama'. We were surprised to find that although 'Drama' is the most popular genre for the Haifa users, they never queried for only 'Drama'.

These findings brought us to the conclusion that we don't have actionable insights to divide the data based on genre preference alone, further research is needed to find if there are sufficient insights by which we can divide the data based on the genres.

(To see the in-depth insights found regarding the genres in the query behaviour and the genres in the movies data – see insight 3 and under Additional Insight > 'Genres – Further analysis on insight 3 for design purposes')

From the additional insight > 'Language' , in the previous question, we can divide horizontally each sites data, based on the language of the movies each sites users query. See 'Additional Insight' > 'Language' output from the previous question.

We chose to divide each site fragments into 2, one containing movies 'English' as one of the spoken languages and the other with all the remaining movies.

We decided to divide by the 'spoken_languages' value, because of finding that 51.77% of users query the data looking only for movies with 'English' as the spoken language. In addition to finding that 36.79% of movies are not available in English. (51.77% was for the overall users, we checked the per location percentage to check if there's any location which wouldn't gain from such fragmentation, found the number varied but was present for every location.)

The result from the additional insight on the fragmentation:

$R_{Tiberias}$:

| Constraints from… | movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Tiberias, Haifa | |
| Insight 2 | | | | | 1990 - current | | | |
| Additional Insight | | | | | | 1. Contains English 2. Other | | |

$R_{Tel-Aviv}$:

| Constraints from… | movie_id | genres | Production _companies | Production_countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Jerusalem, Tel-Aviv | |
| Insight 2 | | | | | 2010 - current | | | |
| Additional Insight | | | | | | 1.Contains English 2. Other | | |

$R_{Jerusalem}$:

| Constraints | movie_id | director |
|---|---|---|
| Movie release date 1990 onwards | | |

| Constraints from… | movie_id | genres | Production _companies | Production_countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Jerusalem, Tel-Aviv | |
| Insight 2 | | | | | 1990 - current | | | |
| Additional Insight | | | | | | 1.Contains English 2. Other | | |

$R_{Eilat}$:

| Constraints from… | movie_id | genres | Production_companies | Production_countries | Release_date | Spoken_languages | cities |
|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Eilat |
| Insight 2 | | | | | 1990 - current | | |
| Additional Insight | | | | | | 1.Contains English 2. Other | |

$R_{Haifa}$:

| Constraints from… | Movie _id | genres | Production _companies | Production _countries | Release _date | Spoken _languages | cities | actors | director |
|---|---|---|---|---|---|---|---|---|---|
| Insight 1 | | | | | | | Tel-Aviv, Haifa, Tiberias | | |
| Insight 2 | | | | | 2010 - current | | | | |
| Additional Insight | | | | | | 1.Contains English 2. Other | | | |

## Part 4 – Final Visualisation

### Tiberias Site Design

**1**

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|

Cities = 'Tiberias' or 'Haifa', **Release_date** $\geq$ 1990,

**Spoken_language** = DOES NOT contain 'English'

**2**

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|

Cities = 'Tiberias' or 'Haifa', **Release_date** $\geq$ 1990,

**Spoken_language** = Contains 'English'

### Tel – Aviv Site Design

**1**

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|

Cities = 'Jerusalem' or 'Tel-Aviv', **Release_date** $\geq$ 2010,

**Spoken_language** = DOES NOT contain 'English'

**2**

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|

Cities = 'Jerusalem' or 'Tel-Aviv', **Release_date** $\geq$ 2010,

**Spoken_language** = Contains 'English'

## Jerusalem Site Design

**1**

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|----------|--------|-----------------------|-----------------------|--------------|------------------|--------|--------|

**Cities** = 'Jerusalem' or 'Tel-Aviv', **Release_date** $\geq$ 1990,

**Spoken_language** = DOES NOT contain 'English'

**2**

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|----------|--------|-----------------------|-----------------------|--------------|------------------|--------|--------|

**Cities** = 'Jerusalem' or 'Tel-Aviv', **Release_date** $\geq$ 1990,

**Spoken_language** = Contains 'English'

**3**

| movie_id | director |
|----------|----------|

Remove all columns except movie_id and director after filtering by: **Release_date** $\geq$ 1990,

## Eilat Site Design

**1**

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities |
|----------|--------|-----------------------|-----------------------|--------------|------------------|--------|

**Cities** = 'Eilat' , **Release_date** $\geq$ 1990,

**Spoken_language** = DOES NOT contain 'English'

**2**

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities |
|----------|--------|-----------------------|-----------------------|--------------|------------------|--------|

**Cities** = 'Eilat' , **Release_date** $\geq$ 1990,

**Spoken_language** = Contains 'English'

## Haifa Site Design

**1**

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors | Director |
|---|---|---|---|---|---|---|---|---|
| | | | | | | | | |

**Cities** = 'Haifa' or 'Tel-Aviv' or 'Tiberias', **Release_date** $\geq$ 2010,

**Spoken_language** = DOES NOT contain 'English'

**2**

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors | Director |
|---|---|---|---|---|---|---|---|---|
| | | | | | | | | |

**Cities** = 'Haifa' or 'Tel-Aviv' or 'Tiberias', **Release_date** $\geq$ 2010,

**Spoken_language** = Contains 'English'

# Appendix 1 – Design prior to genres further research – NOT IMPLEMENTED

## Kibuts Gesher Site Design

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|
| **Cities** = 'Tiberias' or 'Haifa', **Release_date** ≥ 1990, **Genres** = 'Documentary', **Spoken_language** = DOES NOT contain 'English' | | | | | | | |

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|
| **Cities** = 'Tiberias' or 'Haifa', **Release_date** ≥ 1990, **Genres** = 'Documentary', **Spoken_language** = Contains 'English' | | | | | | | |

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|
| **Cities =** 'Tiberias' or 'Haifa', **Release_date** ≥ 1990, **Genres =** 'Family', **Spoken_language** = Contains 'English' | | | | | | | |

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|
| **Cities =** 'Tiberias' or 'Haifa', **Release_date** ≥ 1990, **Genres =** 'Family', **Spoken_language** = DOES NOT contain 'English' | | | | | | | |

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|
| **Cities =** 'Tiberias' or 'Haifa', **Release_date** ≥ 1990, **Genres =** NOT ('Documentary' or 'Family'), **Spoken_language** = Contains 'English' | | | | | | | |

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|
| **Cities =** 'Tiberias' or 'Haifa', **Genres =** NOT ('Documentary' or 'Family'), **Release_date** ≥ 1990, **Spoken_language** =DOES NOT Contain 'English' | | | | | | | |

# Tel – Aviv Site Design

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|
| **Cities** = 'Jerusalem' or 'Tel-Aviv', **Release_date** ≥ 2010, **Genres** = 'Action', **Spoken_language** = DOES NOT contain 'English' | | | | | | | |

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|
| **Cities** = 'Jerusalem' or 'Tel-Aviv', **Release_date** ≥ 2010, **Genres** = 'Action', **Spoken_language** = Contains 'English' | | | | | | | |

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|
| **Cities** = 'Jerusalem' or 'Tel-Aviv', **Release_date** ≥ 2010, **Genres =** NOT 'Action', **Spoken_language** = Contains 'English' | | | | | | | |

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|---|---|---|---|---|---|---|---|
| **Cities** = 'Jerusalem' or 'Tel-Aviv', **Release_date** ≥ 2010, **Genres =** NOT 'Action', **Spoken_language** = DOES NOT contain 'English' | | | | | | | |

## Jerusalem Site Design

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|----------|--------|----------------------|----------------------|--------------|------------------|--------|--------|
| **Cities** = 'Jerusalem' or 'Tel-Aviv', **Release_date** ≥ 1990, **Spoken_language** = DOES NOT contain 'English' | | | | | | | |

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors |
|----------|--------|----------------------|----------------------|--------------|------------------|--------|--------|
| **Cities** = 'Jerusalem' or 'Tel-Aviv', **Release_date** ≥ 1990, **Spoken_language** = Contains 'English' | | | | | | | |

## Eilat Site Design

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities |
|----------|--------|----------------------|----------------------|--------------|------------------|--------|
| **Cities** = 'Eilat' , **Release_date** ≥ 1990, **Spoken_language** = DOES NOT contain 'English' | | | | | | |

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities |
|----------|--------|----------------------|----------------------|--------------|------------------|--------|
| **Cities** = 'Eilat' , **Release_date** ≥ 1990, **Spoken_language** = Contains 'English' | | | | | | |

# *Haifa Site Design*

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors | Director |
|---|---|---|---|---|---|---|---|---|

**Cities** = '  'Haifa'   or 'Tel-Aviv' or 'Tiberias', **Release_date** ≥ 2010,

**Genres** = 'Drama', **Spoken_language** = DOES NOT contain 'English'

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors | Director |
|---|---|---|---|---|---|---|---|---|

**Cities** = '   'Haifa'   or 'Tel-Aviv' or 'Tiberias', **Release_date** ≥ 2010,

**Genres** = 'Drama', **Spoken_language** = Contains 'English'

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors | Director |
|---|---|---|---|---|---|---|---|---|

**Cities** =    'Haifa'    ' or 'Tel-Aviv' or 'Tiberias', **Release_date** ≥ 2010,

**Genres =** NOT 'Drama', **Spoken_language** = Contains 'English'

| movie_id | genres | Production _companies | Production _countries | Release_date | Spoken_languages | cities | actors | Director |
|---|---|---|---|---|---|---|---|---|

**Cities** = '   'Haifa'    ` or 'Tel-Aviv' or 'Tiberias', **Release_date** ≥ 2010,

**Genres =** NOT 'Drama', **Spoken_language** = DOES NOT contain 'English'