# Movies Analysis

**Elnaz Dehkharghani**

## Introduction:

Data is required everywhere to draw meaningful insights about any business. Afterall a useful and insightful analysis of data cannot begin unless the data becomes available. And it is known to all, how important it is to find interesting patterns and visualize the data to gain success in any field. But the flow of data from one place to other is not as simple as it sounds to be. Handling data is treacherous as so many things can go wrong while the transportation. Data can get corrupted, duplicated or accidentally deleted which might increase the latency in order to meet the performance requirements. We all know how difficult it is to work with Big data and imagining a process to extract, transform, aggregate, validate etc. as a manual process can be seen as unrealistic approach.

So here we have a data pipeline, which solves this issue at some level by eliminating many manual steps from the process and enable a smooth and automated flow of data from one place to another. It views all the data as streams of data and give us the flexibility to store the data with any schema.

## Chosen dataset:

We chose movie dataset for building our pipeline. There are two main sources of data.

- API Call
  We have used the TMDB API [1] calls to fetch data for various movies.
- Web Scrapping
  We have scraped the data from the "The Numbers" [2] website.

## Motivation:

Films can be an amazing source of entertaining people due to their key plot messages and unforgettable acting performance in some master pieces. However, they are also a significant part of the business and as such, their main goal is to be profitable. That's why understanding the patterns in successful movies is a crucial step for assisting businesses in establishing the basis of more successful movies. While there are lot of data available on internet about movies provided by IMDB, TMDB, Rotten tomatoes and what not, do we ever wonder what is helping these movies becoming a massive hit? Is there a pattern we can find in the success story of all the movies? We all know that we all have evolved over time, so does our taste in everything, so it might be interesting to see how the love for different kind of movies is changing over time.

Therefore, some of the most interesting questions to answer are:

- What are the most relevant features in successful movies?
- Does budget determine the success of a movie?
- What is the people's satisfaction regarding movies with more budget?
- If the popularity of a movie is helping it earn more profit?

## Data Sources:
- Movies API (discover): Json format result to get information about movies.
- Movies API (movie):  Json format result to get detailed information about a specific movie
- Movies Website (To determine the budget and revenue): Tabular format result to get the budget and revenue of a specific movie

## Tools:
- Python: For calling API requests, manipulate and visualize data
- Beautiful Soup: For web scraping budget and revenue data.
- Apache Kafka: Data Integration tool and data preprocessing
- MongoDB: For storing unstructured results as documents.

## Methodology:
We can divide our project broadly in three different steps.
- Ingesting data from different data sources.
- Streaming data from the different data sources to Mongo DB through Apache Kafka.
- Visualizing the unstructured data stored

## Steps:
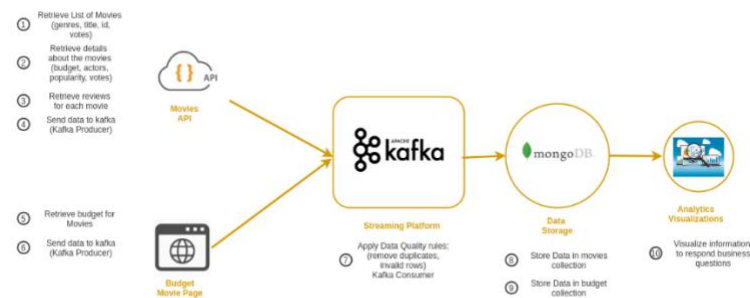The following diagram describe the necessary steps to answer our business questions.



Figure 1.1 General architecture of the system. Steps are indicated for each stage.

- Get Kafka up and running: Lack of knowledge in this tool and complications in Kafka configuration.
- Not a clear understanding in Kafka architecture (consumers and producers) to exchange messages for data processing.
- Kafka connectors configuration to enable data transfer directly to MongoDB.
- Setting up Docker images and modify their configuration for Kafka Confluent.

Decisions:
- Discard Kafka Confluent after having difficulties in installing it and configuring its Mongo Connector as a Docker image. Run the code locally.
- Originally text analysis was considered in our movies' reviews from Reddit, but due to the complexity of the topic and the deadline of the project we decided to put it aside for now.

## Data Ingestion:
We are writing scripts in python to fetch data from two data sources.

### Requesting the Movie API:
In order to receive an API key from the "The MovieDB" website, we needed to create an account. After that, in the API tab or in the provided email address the API Key will be appeared.
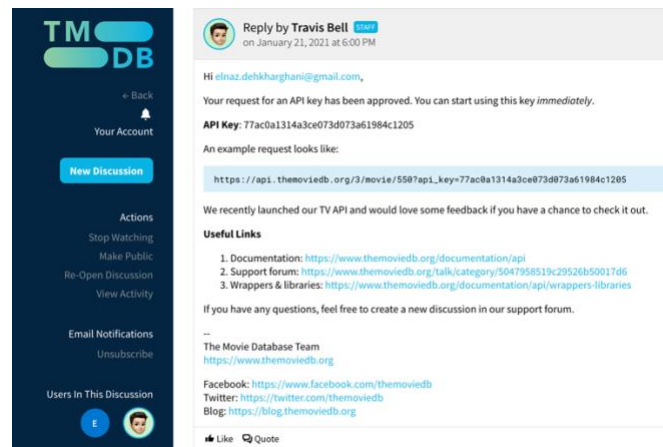


Figure 1.2 Requesting Movies API key to access data.

First, one deals with calling the TMDB (the movie data base) API. We are basically calling two of its API requests. Calling both the APIs require an API key we already received in the previous step.

- **GET**/discover/movie

    Discover API discovers movies by different types of data. It gives an overview about title of the movie, release date, popularity, vote count. We can also apply filters for release data, title, page value etc. to narrow down the results. An example of the API call can be seen in the below figure.
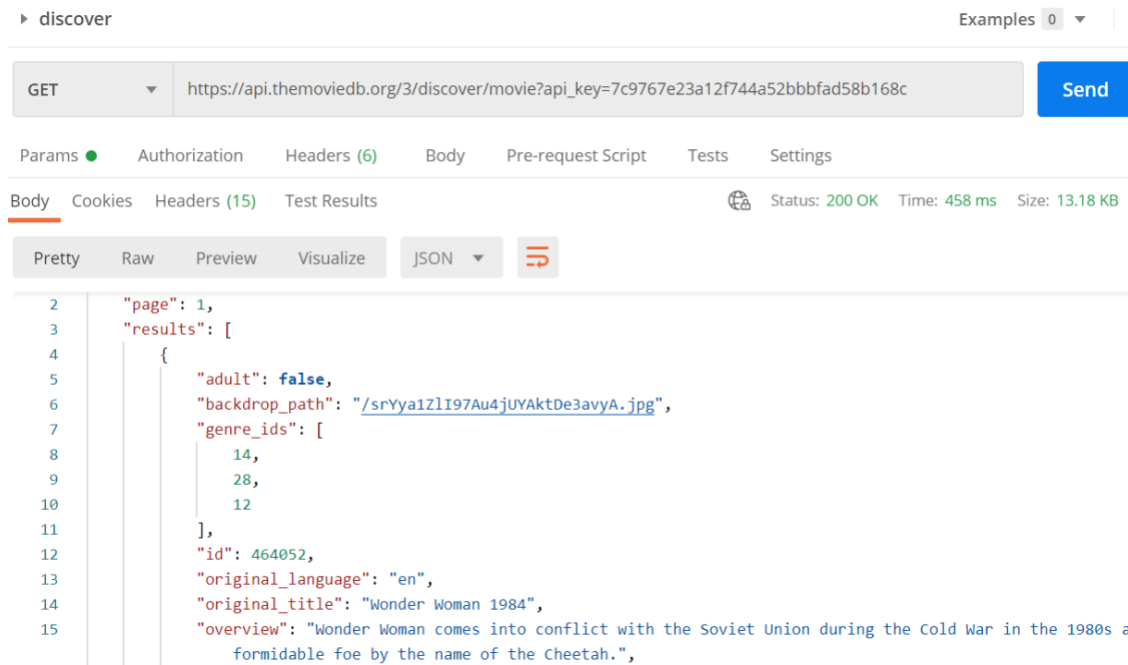
Figure 1.3 Example of the /discover/movie API call.

We have stored the data from the year 2000 till 2021. Using this API, we are storing the ids of all the movies that were released in this duration as the id field is the unique identifier here which we can use in the next API call.

Only the unique ids are being stored to avoid duplicates in data.

```python
# ====== STEP 1 : Get List of movies ============
print('Step 1: Loading movies list')
# Get list of moviews

for i in range(2000, 2021, 1):
    release_year = i
    sort_by = 'vote_average.desc'
    url_search = f'https://api.themoviedb.org/3/discover/movie?primary_release_year={release_year}&sort_by={sort_by}&api_key={config.

    # Call the API
    print(url_search)
    res = requests.get(url_search)
    json_obj = res.json()
    results = json_obj['results']

    # Add only unique movies id
    for result in results:
        if result['id'] not in movies:
            movies.append(result['id'])
```

Figure 1.4 Remove duplicated movies ids after calling /discover/movie API.

- **GET**/movie/{movie_id}

After we have all the ids of the movies released in the year 2000 till 2021, we can call the above API with the movie id. It will give a more detailed response for the movies which will include information like budget, production companies, production countries, revenue etc.

There is also an optional query parameter `append_to_response` which helps us get information like cast and crew members who took part in that movie.

```python
movies = []
for i in range(page, total_pages):
    print(f'Page: {i}')
    # Call >= 2 page
    if (i > 1):
        url_search = f'https://api.themoviedb.org/3/discover/movie?primary_release_year={release_year}&sort_by={sort_by}&api_key={config.API
        # Add page to get results
        url_search += f'&page={i}'
        res = requests.get(url_search)
        json_obj = res.json()
        results = json_obj['results']

    # Add only unique movies id
    for result in results:
        if result['id'] not in movies:
            movies.append(result['id'])

# Print results
print(f'Total movies added: {len(movies)}')
for m in movies:
    print(m)
```

Figure 1.5 Obtain movie details by calling /movie/movie_id API.

Results of both the APIs were stored as a JSON object. Also, we thought of just keeping the useful fields for sending data to Kafka for visualization and not sending the entire response. The fields we sent to Kafka were as follows:
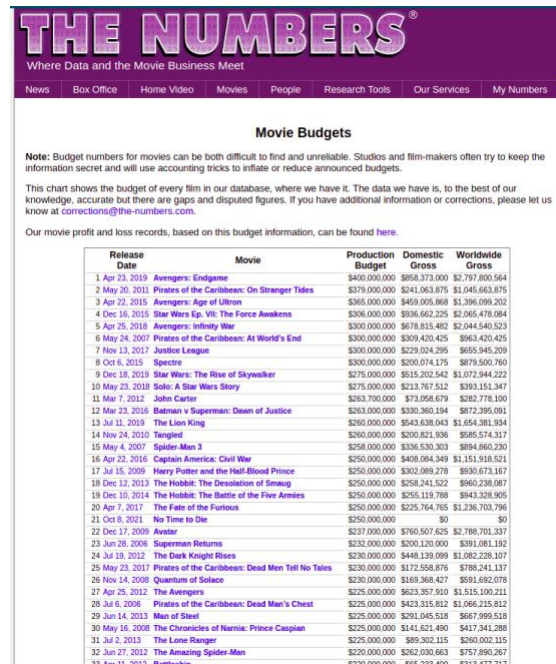
```python
# movies_list = []
for m_id in movies:
    url_search = f'https://api.themoviedb.org/3/movie/{m_id}?api_key={config.API_KEY}&append_to_response=credits'
    res = requests.get(url_search)
    result = res.json()
    movie = {}
    try:
        movie['_id'] = str(result['id'])
        movie['budget'] = result['budget']
        movie['genres'] = list(map(lambda m: m['name'], result['genres']))
        movie['original_language'] = result['original_language']
        movie['original_title'] = result['original_title']
        movie['popularity'] = result['popularity']
        movie['production_companies'] = list(map(lambda m: m['name'], result['production_companies']))
        movie['production_countries'] = list(map(lambda m: m['name'], result['production_countries']))
        movie['release_date'] = result['release_date']
        movie['revenue'] = result['revenue']
        movie['runtime'] = result['runtime']
        movie['status'] = result['status']
        movie['vote_average'] = result['vote_average']
        movie['vote_count'] = result['vote_count']
        movie['cast'] = list(map(lambda m: {'department': m['known_for_department'], 'name': m['original_name'],
                                            'popularity': m['popularity']}, result['credits']['cast']))
        movie['crew'] = list(map(
            lambda m: {'department': m['department'], 'name': m['original_name'], 'popularity': m['popularity']},
            result['credits']['crew']))
        message = json.dumps(movie)
        print(f" ======== \n {movie['original_title']}")
```

Figure 1.6 Fields of interest of each movie to send to Kafka.

## Budget Movies

The data we received from the API had a lot of null values for the field budget and revenue. These fields were important parts of our visualization analysis. Though removing all those documents having null value for budget and revenue would have incurred a significant data loss.

We tried to gather data from a website for the missing values of budget and revenue. The source was "The Numbers where data and the movie business meet" [2]. We needed to do web data scraping to fulfil our requirement. The advantage is that the web page structure makes this process achievable, and the information can be easily obtained. A screenshot of the webpage is shown below:



Figure 1.7 List of movies from the webpage to scrap

The general steps for this web scrapping are listed as follow:
1. Use splinter and beautiful soup python modules inside a Kafka producer to access the html content of a single webpage. The movies are paginated like 101, 201, 301, etc.
2. For each page, get Title, Production budget, Domestic Gross and Worldwide Gross.
3. Send data to Kafka and process the information.

## Data Source Automation:

We have automated the data sourcing process by using crontab library of python. The example of which can be seen in the underlined script. These scripts will run once every day at the second hour.



Figure 1.8 Automation process for movies producer

# Data Storage:

## Mongo DB

As the majority of the information comes in JSON format, it is much easier to handle the storage if we use an unstructured data base like Mongo DB. In addition, the whole team wanted to learn this technology since it is more appealing to work in Big Data use cases.

Specifically, we used a python module called pymongo which allowed us to interact directly with the database from the Kafka consumers, reducing the difficulties with the connection between Kafka and Mongo DB rather than using a Kafka connector (as indicated in our difficulties).

To work simultaneously in the same database, we created a cluster in Mongo DB Atlas where the team could insert and retrieve information independently and develop different parts of the project. As we have two sources of information (movies API and budget movies) we defined 2 data collections: movies_collection and budget_collection.

**Movies_collection** => detailed movies information (4,729 movies as of Feb 7, 2021)
**Budget_collection** => financial movies information (4,890 movies as of Feb 7, 2021)

A screenshot of the collections is shown below:



Figure 1.9 Movies Collections in Mogo Atlas. Left Image - Detailed data related to movies.
Right Image – Financial movies data

## Kafka producers and consumers

Kafka was selected to successfully handle the data processing. As the Kafka Architecture suggests, we need to use Producers and Consumers to exchange messages of different topics, handle them and target them to their respective data storage. The following diagram shows how the producers and consumers were related:

Figure 1.10 Kafka producers and consumers.

## Kafka Configuration:

As mentioned before, we decided to work with Kafka, for this purpose we tried 2 approaches to setting up the environment: Locally and in Docker. But as we continued encountering difficulties in setting the MongoDB connector for Confluent, we decided to run it locally.

## Local Environment:

Initially Kafka is downloaded from the mirror site. After downloading and unzipping the Kafka through the terminal you must go to the Kafka folder to run the zookeeper and the Kafka broker service to start all services in the correct order:



Figure 1.11 Kafka services running: Left Image – Zookeeper service. Middle Image – Broker service. Right Image – New topic created. All commands are in independent terminals.

## Testing the Producer and Consumer:

A Kafka client communicates with the Kafka brokers via the network for writing (or reading) events. Once received, the brokers will store the events in a durable and fault-tolerant manner for as long as you need—even forever. Run the console producer client to write a few events into your topic. By default, each line you enter will result in a separate event being written to the topic.

Figure 1.12 Example of sending and receiving a message in Kafka. Left Image – Producer. Right Image - Consumer

Now we do the same process by our movie API:



Figure 1.13 Example of sending movies messages in Kafka.

After we read movies information in the Kafka consumer, we send the data directly to the Kafka Broker where the consumer performed a simple data pre-processing, including a validation in our DB to avoid duplicates.



Image 1.14. Budget Consumer in python

## Kafka Consumer (movies)

In order to store the "movies" topic information, a python consumer listens to all incoming messages from the producer and inserts a batch of movies into Mongo DB after validating whether that movie already exists or not. Part of the code is shown in the image below:

## Kafka Consumer (budget)

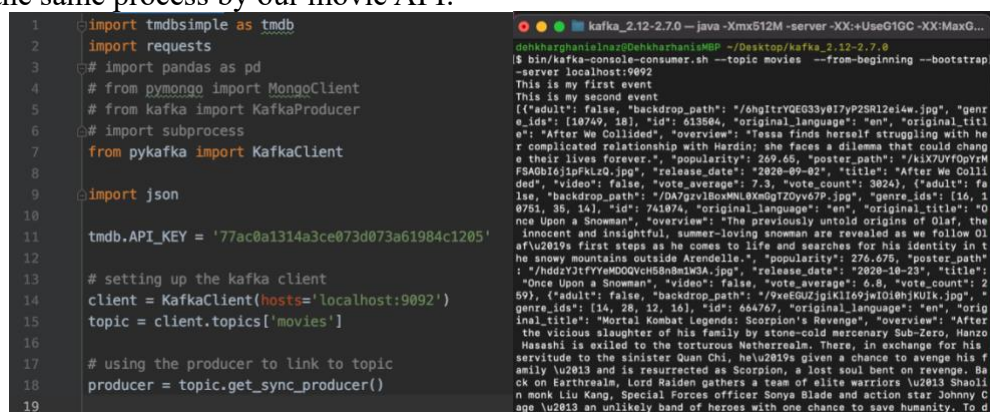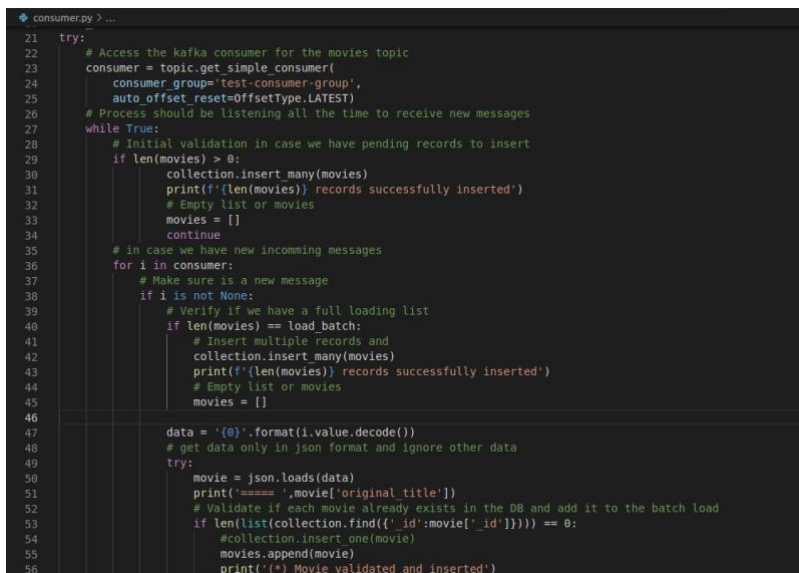In order to store the "budget" topic information, we wrote a python consumer which listens to all incoming messages from the producer and inserts a batch of movies with budgets into Mongo after a validation to determine whether that movie already exists or not. Part of the code is shown in the image below:



Figure 1.15 Kafka consumer for the budget producer

## Aggregated Collection

Since we have information about movies and financial results of some of them, it is possible to merge that information in another collection. This is mainly because not all the movies results appear in the from the financial movies collection. That is why we decided to create a new collection called aggregated_collection to store the data from the movies_collection joint with the budget_collection.



Figure 1.16 Aggregated collection that combines data from movies and budget collection.

The python code that creates this collection is shown below: (We have applied the left join on the two collections, movies and budget_collection)

```
client = pymongo.MongoClient(conn)
db = client["movies"]
collection1 = db.movies_collection
collection2 = db.budget_collection
collection3 = db.aggregated_collection


data = list(collection1.aggregate([

    {"$lookup": {
        "from": "budget_collection",
        "localField": "original_title",
        "foreignField": "title",
        "as": "aggregate"
    }

    },
]))


collection3.insert_many(data)
print(f'{len(data)} is the number of movies inserted in the aggregated_collection')
```

Figure 1.17 Python code to create a new aggregated collection.

## Access data in Mongo

One step before visualizing data is to define the approach of getting the data we need from Mongo. As we are using python as our main tool to manipulate data, we can build a custom API in which we define the functions we need to get the data in an appropriate format that pandas can process.

Therefore, the following functions retrieve results for only movies in English:

| Custom API method | Explanation |
|---|---|
| getMoviesReleaseYear(topN=10) | Retrieves the top 10 movies based on their released year and highest vote average. |
| getTopGenres(topN=10) | Retrieves the top 10 genres based on the highest number of movies and highest average revenue. |
| getTopProductionCompanies(topN=10) | Retrieves the top 10 production companies based on the highest number of movies released and highest average revenue. |
| getTopMoreRevenue(topN=10) | Retrieves the top 10 movies based on their release year and highest revenue. |
| getTopVoteAvg(topN=10) | Retrieves top 10 movies based on their release year and highest vote average. |
| wholeData(topN=1000) | Retrieves the top 1,000 movies based on their release date and vote average. |

An example of the python code for some of these functions is shown in the following picture:

```
49  # Function to get top 10 genres
50  def getTopGenres(topN=10):
51      movies = list(collection.aggregate([{"$unwind": "$genres"}, {"$group": {"_id": "$genres",
52                      "avgRuntime": {"$avg": "$runtime"},
53                      "noMovies": {"$sum": 1},
54                      "avgPopularity": {"$avg": "$popularity"},
55                      "avgVotesNo": {"$avg": "$vote_count"},
56                      "avgProductionBudget": {"$avg": {"$arrayElemAt": ['$aggregate.productionBudget', 0]}},
57                      "avgDomesticBudget": {"$avg": {"$arrayElemAt": ['$aggregate.domesticBudget', 0]}},
58                      "avgWorldwideGross": {"$avg": {"$arrayElemAt": ['$aggregate.worldwideGross', 0]}},
59                      "avgBudget": {"$avg": "$budget"},
60                      "avgRevenue": {"$avg": "$revenue"},
61                      "avgVotes": {"$avg": "$vote_average"}}},
62                      {"$sort": {"noMovies": -1, "avgRevenue": -1}}]))
63      movies_json = []
64      # In case we get less results than the especified number
65      topN = (len(movies)) if len(movies) < topN else topN
66      for i in range(topN):
67          movie_json = {
68              "genres": movies[i]["_id"],
69              "avgRuntime": round(movies[i]["avgRuntime"], 3),
70              "noMovies": movies[i]["noMovies"],
71              "avgPopularity": round(movies[i]["avgPopularity"], 3),
72              "avgVotes": round(movies[i]["avgVotesNo"], 3),
73              "avgBudget": round(movies[i]["avgBudget"], 3),
74              # new fields from aggregation
75              "avgProductionBudget": (0 if movies[i]["avgProductionBudget"] == None else round(movies[i]["avgProductionBudget"],3)),
76              "avgDomesticBudget": (0 if movies[i]["avgDomesticBudget"] == None else round(movies[i]["avgDomesticBudget"],3)),
77              "avgWorldwideGross": (0 if movies[i]["avgWorldwideGross"] == None else round(movies[i]["avgWorldwideGross"],3)),
78              "avgRevenue": round(movies[i]["avgRevenue"], 3),
79              "avgVotes": round(movies[i]["avgVotes"], 3)
80          }
81          movies_json.append(movie_json)
82      return movies_json
```

Figure 1.18 Python code of getTopGenres() function. We use pymongo to access the data in our Mongo DB and perform some aggregations when needed

## Data Visualization

In the upcoming stages we want to process our data to find the answers to the following questions:

1. Which movie has the highest and lowest profit, budget, revenue, and runtime in our dataset?
2. What are the top 10 movies in profit, budget, revenue, and runtime?
3. What are the top 30 movies in popularity?
4. What are the top 30 movies in vote_count?
6. What is the average budget, profit, popularity, etc for top 10 Genres?
7. Comparing the average Budget, Domestic Gross, Worldwide Gross for top Genre in the single plot?
8. What is the average budget, profit, popularity, etc for top 10 Companies?
9. Comparing the average Budget, Domestic Gross, Worldwide Gross for top companies in the single plot?
10. What are the key features of the movie success rate for our future modeling?

Until now, with the help of Kafka we could ingest our data from two different sources, calling API and web scraping and then store them in the Atlas MongoDB. Therefore, it is time to visualize our data to have a good insight into our collected data and find the answers to the above questions. For this purpose, we need NumPy, Pandas, Seaborn, and Matplotlib libraries. Beforehand to avoid any incorrect result in our data analysis we need to do one of the most important tasks: data cleaning. This step usually takes too much time and effort of any data scientist and most of the times it takes more than 70 percent of his time.

```
import math
from datetime import datetime
import json
from pymongo import MongoClient
import dbQueries
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from matplotlib import gridspec
from collections import Counter
sns.set()
```

```
plt.style.use('seaborn-notebook')
```

Figure 1.19 Elements used in the project. Left Image – python modules. Right Image – Style for plots.

First Part (Visualization – Investigating on Movies):

Loading Whole Data from Mongodb:

To have access to our MongoDB we need to write functions in a separate file called dbQueries.py which is an API to make movies data available to anyone who wants to use it. We mainly used pymongo module to make a client in order to fetch our data directly from MongoDB:

```
def wholeData(topN=1000):
    # Filter only movies in english
    movies = list(collection.find({"original_language":"en"}))
    movies_json = []
    # In case we get less results than the especified number
    topN = (len(movies) if len(movies) < topN else topN)
    for i in range(topN):
        movie_json = {
            "title": movies[i]["original_title"],
            "genres": movies[i]["genres"],
            "budget": movies[i]["budget"],
            "popularity": movies[i]["popularity"],
            "production_companies": movies[i]["production_companies"],
            "production_countries": movies[i]["production_countries"],
            "revenue": movies[i]["revenue"],
            # New fields from aggregated collection
            "productionBudget": (0 if len(movies[i]["aggregate"]) == 0 else round(movies[i]["aggregate"][0]['productionBudget'],3)),
            "domesticBudget": (0 if len(movies[i]["aggregate"])== 0 else round(movies[i]["aggregate"][0]['domesticBudget'],3)),
            "worldwideGross": (0 if len(movies[i]["aggregate"]) == 0 else round(movies[i]["aggregate"][0]['worldwideGross'],3)),
            "runtime": movies[i]["runtime"],
            "status": movies[i]["status"],
            "vote_avg": movies[i]["vote_average"],
            "vote_count": movies[i]["vote_count"],
            "release_date": movies[i]["release_date"]
        }
        movies_json.append(movie_json)
    return movies_json
```

Figure 1.20 - Whole Data API code implementation.

The wholeData function has access to 1,000 of our aggregated data (aggregated_collection). These data are aggregated from two other collections by "Left Join" of (budget_collection and movies_collection) in the Atlas MongoDB. For our analysis, we are only considering movies in English language.

movies

aggregated_collection

budget_collection

movies_collection

Figure 1.21 - MongoDb collections

Before going ahead, we need to know our data set:

**title:** The Official Title of the movie.
**genres:** A stringified list of dictionaries that list out all the genres associated with the movie.
**budget:** The budget of the movie in dollars.
**popularity:** The Popularity Score assigned by TMDB.
**production_companies:** A stringified list of production companies involved with the making of the movie.
**production_countries:** A stringified list of countries where the movie was shot/produced in.
**revenue:** The total revenue of the movie in dollars.

**Production budget:** The budget of the movie in dollars. (Added by web scraping from another source)
**Domestic Gross:** The total revenue of the movie in dollars in domestic area. (Added by web scraping from another source)
**Worldwide Gross:** The total revenue of the movie in dollars in worldwide. (Added by web scraping from another source)
**runtime:** The runtime of the movie in minutes.
**status:** The status of the movie (Released, To Be Released, Announced, etc.)
**vote_average:** The average rating of the movie.
**release_date:** Theatrical Release Date of the movie.

Now we import and call this wholedata function in our Jupyter notebook for further analysis, and then convert our JSON data to a data Frame object. In case to have an initial insight and make sure that our data is loaded properly we use the head() function.

```python
re_json = dbQueries.wholeData()
wholeData = pd.DataFrame.from_dict(re_json)
wholeData.head(5)
```

Figure 1.22 - Transform API result into dataframe in python

Before doing any data cleaning we would have a quick look at the number of rows, columns, column names, and their data types.

```python
wholeData.shape
wholeData.columns
wholeData.dtypes
```

Figure 1.23 - python commands to understand the data structure

To check how many rows have null values we use this function.

```python
wholeData.isnull().sum()
```

## Data cleaning:

Changing the data type of the release date column to datetime format.

```python
wholeData['release_date'] = pd.to_datetime(wholeData['release_date'])
```

Since some values in the budget column are missing, we can fill them with the productionBudget column that is fetched from the web scraping source.

```python
wholeData['budget'] = np.where(wholeData['budget'] == 0, wholeData['productionBudget'], wholeData['budget'])
```

After that, if still there are some missing values in the rows, we replace those 0 values in both revenue and budget columns with NaN values and then drop them.

```python
wholeData[['budget','revenue']] = wholeData[['budget','revenue']].replace(0,np.NAN)

wholeData.dropna(subset=['budget', 'revenue'], inplace=True)
print('After cleaning, we have {} rows'.format(wholeData.shape[0]))
```

In this step for further analysis, we need to add a new column as Profit for each of the movies. We do this by subtracting the revenue column by budget column. Also, to change the datatype of some columns to int64 we apply the Numpy.int64 function.

```
wholeData['profit'] = wholeData['revenue']-wholeData['budget']
wholeData['profit'] = wholeData['profit'].apply(np.int64)
wholeData['budget'] = wholeData['budget'].apply(np.int64)
wholeData['revenue'] = wholeData['revenue'].apply(np.int64)
wholeData['productionBudget'] = wholeData['productionBudget'].apply(np.int64)
wholeData['domesticBudget'] = wholeData['domesticBudget'].apply(np.int64)
wholeData['worldwideGross'] = wholeData['worldwideGross'].apply(np.int64)
```

## Data Analysis:

By using the describe() function, we have a quick summary of statistics related to the DataFrame columns. This function gives the mean, std and IQR values.

Before going ahead, we need to define two different functions: one for finding the minimum and maximum values of a given column and the other for finding the first top ten for that specific column name. In the first function by idxmin() and idxmax() we get the id of the minimum and maximum values in the column. After that by loc[] we can filter our data frame to catch only these two rows by all their columns then we concatenate them together in a single data frame. For the second function we need to sort the values in the descending order and just define a size as given parameter to catch results. Also, with Numpy we applied the mean function to demonstrate where the mean value is in the x-axis. At the end with seaborn we will plot the result with that previously specified figure size.

```python
def find_min_max(col_name):
    min_index = wholeData[col_name].idxmin()
    max_index = wholeData[col_name].idxmax()
    low  = pd.DataFrame(wholeData.loc[min_index,:])
    high = pd.DataFrame(wholeData.loc[max_index,:])
    print('Movie which has highest '+col_name+' : ', wholeData['title'][max_index])
    print('Movie which has lowest '+col_name+' : ', wholeData['title'][min_index])
    return pd.concat([high,low], axis=1)
```

```python
def top_10(col_name,size=10):
    df_sorted = pd.DataFrame(wholeData[col_name].sort_values(ascending=False))[:size]
    df_sorted['title'] = wholeData['title']
    plt.figure(figsize=(12,6))
    #Calculate the avarage
    avg = np.mean(wholeData[col_name])
    sns.barplot(x=col_name, y='title', data=df_sorted, label=col_name)
    plt.axvline(avg, color='k', linestyle='--', label='mean')
    if (col_name == 'profit' or col_name == 'budget' or col_name == 'revenue'):
        plt.xlabel(col_name.capitalize() + ' (U.S Dolar)')
    else:
        plt.xlabel(col_name.capitalize())
    plt.ylabel('')
    plt.title('Top 10 Movies in: ' + col_name.capitalize())
    plt.legend()
```

We will apply these two functions on the following metrics: profit, budget, revenue, and runtime and the illustrated results will be as below. Since we are fetching data in a real time streaming by Kafka messenger through MovieDB API, results will be changed at any time that the program is running.
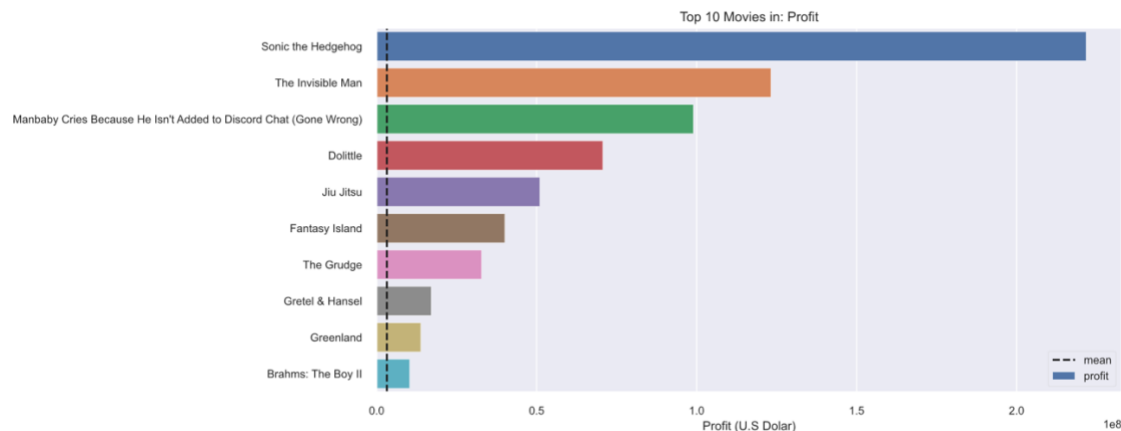
Results for the movies with highest and lowest Profit:

```
find_min_max('profit')

Movie which has highest profit :  Sonic the Hedgehog
Movie which has lowest profit :  Mulan
```

| | 64 | 17 |
|---|---|---|
| title | Sonic the Hedgehog | Mulan |
| genres | [Action, Science Fiction, Comedy, Family] | [Adventure, Fantasy] |
| budget | 85000000 | 200000000 |
| popularity | 125.925 | 332.695 |
| production_companies | [Original Film, Blur Studios, Marza Animation ... | [Walt Disney Pictures, China Film Group Corpor... |
| production_countries | [Japan, United States of America] | [China, United States of America] |
| revenue | 306766470 | 57000000 |
| productionBudget | 0 | 200000000 |
| domesticBudget | 0 | 0 |
| worldwideGross | 0 | 69963008 |
| runtime | 99.0 | 115.0 |
| status | Released | Released |
| vote_avg | 7.5 | 7.1 |

## Results for the top 10 movies in Profit:



## Results for the movies with highest and lowest Budget:

```
find_min_max('budget')

Movie which has highest budget :  Mulan
Movie which has lowest budget :  Manbaby Cries Because He Isn't Added to Discord Chat (Gone Wrong)
```

| | 17 | 381 |
|---|---|---|
| title | Mulan | Manbaby Cries Because He Isn't Added to Discor... |
| genres | [Adventure, Fantasy] | [War] |
| budget | 200000000 | 1 |
| popularity | 332.695 | 0.6 |
| production_companies | [Walt Disney Pictures, China Film Group Corpor... | [Alwharf Studios] |
| production_countries | [China, United States of America] | [] |
| revenue | 57000000 | 99000000 |
| productionBudget | 200000000 | 0 |
| domesticBudget | 0 | 0 |
| worldwideGross | 69963008 | 0 |
| runtime | 115.0 | 34.0 |
| status | Released | Released |
| vote_avg | 7.1 | 10.0 |

Results for the top 10 movies in Budget:


Top 10 Movies in: Budget

Results for the movies with highest and lowest Revenue:



```
find_min_max('revenue')

Movie which has highest revenue :  Sonic the Hedgehog
Movie which has lowest revenue :  Rocky Horror Remade
```

|  | 64 | 407 |
| --- | --- | --- |
| title | Sonic the Hedgehog | Rocky Horror Remade |
| genres | [Action, Science Fiction, Comedy, Family] | [Horror, Animation, Music, Comedy, Science Fic... |
| budget | 85000000 | 1 |
| popularity | 125.925 | 1.96 |
| production_companies | [Original Film, Blur Studios, Marza Animation ... | [The Player's Society] |
| production_countries | [Japan, United States of America] | [United States of America] |
| revenue | 306766470 | 1 |
| productionBudget | 0 | 0 |
| domesticBudget | 0 | 0 |
| worldwideGross | 0 | 0 |
| runtime | 99.0 | 95.0 |
| status | Released | Released |
| vote_avg | 7.5 | 10.0 |

Results for the top 10 movies in Revenue:



Results for the movies with highest and lowest Runtime:

Results for the top 10 movies in Runtime:



Results for the first top 30 movies in Popularity:



Results for the first top 30 movies in vote_count:

For plotting the pairwise relationships between given columns in a dataset, we have applied the seaborn.paiplot() function and gave a parameter as kind: reg to have regression plots.

```python
df_related = wholeData[['profit','budget','revenue','runtime', 'vote_count','popularity']]

sns.pairplot(df_related, kind='reg')
#plt.savefig('correlation.png',bbox_inches='tight')
```



To visualize only one of these plots here for budget, revenue in a single figure we used this regplot() function of seaborn.

```python
sns.regplot(x=wholeData['revenue'], y=wholeData['budget'],color='c')
```

## Results for key features of success rate a movie:

We defined a function to plot a heatmap by finding the correlation between given variables here as our columns of the data frame. A correlation plot used to investigate the dependency between multiple variables at the same time and to highlight the most correlated variables in a data frame.

*# sns.diverging_palette:* Make a diverging palette between two HUSL colors.

*# as_cmapbool:* If True, return a matplotlib.colors.Colormap.

*# Sns.heatmap:* Plot rectangular data as a color-encoded matrix.

*# cbar:* Whether to draw a colorbar.

*# ax:* Axes in which to draw the plot, otherwise use the currently-active Axes

*# annot:* If True, write the data value in each cell.

*# annot_kws:* Keyword arguments for matplotlib.axes.Axes.text() when annot is True.

*# Square:* If True, set the Axes aspect to "equal" so each cell will be square-shaped.

```python
def plot_correlation_map( wholeData ):
    corr = wholeData.corr()
    _ , ax = plt.subplots( figsize =( 12 , 10 ) )
    cmap = sns.diverging_palette( 240 , 10 , as_cmap = True )
    _ = sns.heatmap(corr,cmap = cmap,square=True, cbar_kws={ 'shrink' : .9 }, ax=ax,
                    annot = True, annot_kws = { 'fontsize' : 12 })
```

```python
plot_correlation_map(wholeData[['popularity','budget','revenue','runtime','vote_count']])
```

Interpretation of this plot:

If we have a Negative correlation it means that features increase or decrease in values in opposite directions on a number line. If one feature increases, the other feature decreases. But on the other hand, Positive correlation means features decrease or increase in values together. If two features have a correlation of one, then the features are duplicates.

As an example, since this plot is showing budget and revenue or budget and vote counts have a positive high correlation, it is obvious that with more investment would have more revenues and vote counts. A higher correlation of features with the target feature indicates the possible higher importance of the feature for any future modeling.

## Second Part: (Visualization – Investigating on Genres):

### Loading Only Ten Top Genre From MongoDB:

Our API can return the top ten genres from our MongoDB collection. The method works as follows: Since the genres for each movie are stored as an array first, we need to apply $unwind to deconstruct them. After that with $group we can group these genres by their ids. For the other values in the document, we calculate their averages and finally sort them base on the number of movies and average revenues, we iterate until to find these ten genres and the result will be saved in a list:

```python
# Function to get top 10 genres
def getTopGenres(topN=10):
    movies = list(collection.aggregate([{"$unwind": "$genres"}, {"$group": {"_id": "$genres",
                "avgRuntime": {"$avg": "$runtime"},
                "noMovies": {"$sum": 1},
                "avgPopularity": {"$avg": "$popularity"},
                "avgVotesNo": {"$avg": "$vote_count"},
                "avgProductionBudget": {"$avg": {"$arrayElemAt": ['$aggregate.productionBudget', 0]}},
                "avgDomesticBudget": {"$avg": {"$arrayElemAt": ['$aggregate.domesticBudget', 0]}},
                "avgWorldwideGross": {"$avg": {"$arrayElemAt": ['$aggregate.worldwideGross', 0]}},
                "avgBudget": {"$avg": "$budget"},
                "avgRevenue": {"$avg": "$revenue"},
                "avgVotes": {"$avg": "$vote_average"}}},
                {"$sort": {"noMovies": -1, "avgRevenue": -1}}]))
    movies_json = []
    # In case we get less results than the especified number
    topN = (len(movies) if len(movies) < topN else topN)
    for i in range(topN):
        movie_json = {
            "genres": movies[i]["_id"],
            "avgRuntime": round(movies[i]["avgRuntime"], 3),
            "noMovies": movies[i]["noMovies"],
            "avgPopularity": round(movies[i]["avgPopularity"], 3),
            "avgVotes": round(movies[i]["avgVotesNo"], 3),
            "avgBudget": round(movies[i]["avgBudget"], 3),
            # new fields from aggregation
            "avgProductionBudget": (0 if movies[i]["avgProductionBudget"] == None else round(movies[i]["avgProductionBudget"],3)),
            "avgDomesticBudget": (0 if movies[i]["avgDomesticBudget"] == None else round(movies[i]["avgDomesticBudget"],3)),
            "avgWorldwideGross": (0 if movies[i]["avgWorldwideGross"] == None else round(movies[i]["avgWorldwideGross"],3)),
            "avgRevenue": round(movies[i]["avgRevenue"], 3),
            "avgVotes": round(movies[i]["avgVotes"], 3)
        }
        movies_json.append(movie_json)
    return movies_json
```

Now with calling this function and applying the head() method we would be sure of our data loading:

```
# Display results per genre
re_json = dbQueries.getTopGenres(10)
df_topGenres = pd.DataFrame.from_dict(re_json)
# ======= Columns added: avgProductionBudget, avgDomesticBudget, avgWorldwideGross =======
df_topGenres.head()
```

| | genres | avgRuntime | noMovies | avgPopularity | avgVotes | avgBudget | avgProductionBudget | avgDomesticBudget | avgWorldwideGross | avgRevenue |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Drama | 81.121 | 1169 | 11.676 | 8.638 | 1574308.010 | 3.377920e+07 | 4.037743e+07 | 7.649464e+07 | 3911999.482 |
| 1 | Documentary | 110.208 | 1054 | 2.223 | 9.383 | 4910.930 | 2.550000e+07 | 3.190621e+07 | 6.845040e+07 | 22916.507 |
| 2 | Comedy | 67.746 | 760 | 17.191 | 8.658 | 2337151.226 | 3.714242e+07 | 4.488247e+07 | 8.158763e+07 | 4024125.428 |
| 3 | Music | 70.907 | 529 | 6.365 | 9.453 | 549995.849 | 5.050000e+07 | 3.263855e+07 | 6.939632e+07 | 1135669.070 |
| 4 | Romance | 88.394 | 325 | 13.211 | 8.599 | 830369.923 | 2.663846e+07 | 3.491543e+07 | 6.057785e+07 | 2699984.206 |

We see here number of columns and rows:

```
df_topGenres.shape

(10, 10)
```

With describe() function find a summary of statistics.

```
df_topGenres.describe()
```

| | avgRuntime | noMovies | avgPopularity | avgVotes | avgBudget | avgProductionBudget | avgDomesticBudget | avgWorldwideGross | avgRevenue |
|---|---|---|---|---|---|---|---|---|---|
| count | 10.000000 | 10.000000 | 10.000000 | 10.00000 | 1.000000e+01 | 1.000000e+01 | 1.000000e+01 | 1.000000e+01 | 1.000000e+01 |
| mean | 73.966000 | 525.400000 | 21.165000 | 8.69860 | 2.355983e+06 | 4.000419e+07 | 3.890693e+07 | 7.783373e+07 | 4.642771e+06 |
| std | 17.432723 | 347.370823 | 12.459941 | 0.41898 | 2.170878e+06 | 1.319696e+07 | 7.055758e+06 | 1.498427e+07 | 4.385676e+06 |
| min | 45.034000 | 237.000000 | 2.223000 | 8.11500 | 4.910930e+03 | 2.550000e+07 | 3.054113e+07 | 6.057785e+07 | 2.291651e+04 |
| 25% | 64.212000 | 299.750000 | 12.059750 | 8.55925 | 9.190526e+05 | 3.036558e+07 | 3.320777e+07 | 6.868688e+07 | 2.243539e+06 |
| 50% | 72.601500 | 321.000000 | 23.342000 | 8.62650 | 1.868793e+06 | 3.546081e+07 | 3.737013e+07 | 7.313044e+07 | 3.899978e+06 |
| 75% | 80.117250 | 702.250000 | 30.569750 | 8.67075 | 2.509509e+06 | 4.879688e+07 | 4.375621e+07 | 8.147007e+07 | 4.491439e+06 |
| max | 110.208000 | 1169.000000 | 37.374000 | 9.45300 | 6.375323e+06 | 6.533333e+07 | 5.037111e+07 | 1.128220e+08 | 1.511300e+07 |

Calculate number of null values:

```
# Checking for null values
df_topGenres.isnull().sum()

genres                 0
avgRuntime             0
noMovies               0
avgPopularity          0
avgVotes               0
avgBudget              0
avgProductionBudget    0
avgDomesticBudget      0
avgWorldwideGross      0
avgRevenue             0
dtype: int64
```

Checking for any duplicate existence:

```
# Checking for duplications
sum(df_topGenres.duplicated())
#df.drop_duplicates(inplace=True)

0
```

## Data Cleaning:

Data cleaning and imputing the average profit column in the data frame would be quite same with which we had before on our previous whole data therefor to avoid any repetition we did not write this part here.

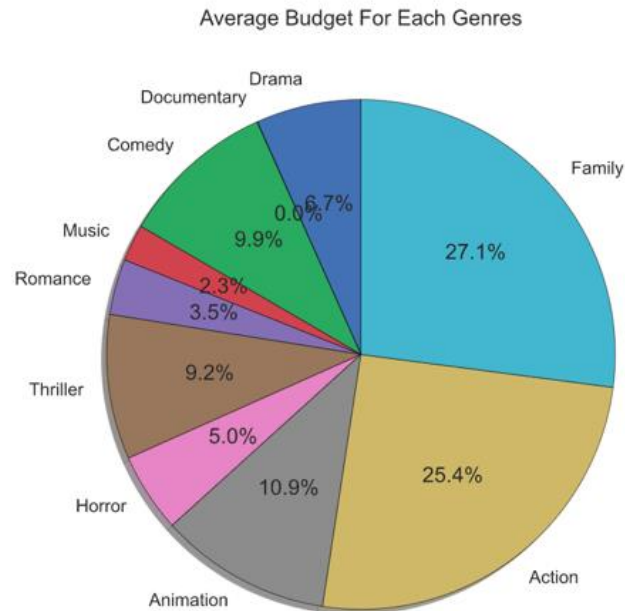Plotting the Average budget for these 10 top Genre with a piechart.

```
slices = df_topGenres['avgBudget']
labels = df_topGenres['genres']

plt.pie(slices, labels = labels,shadow= True, startangle=90, autopct='%1.1f%%', wedgeprops={'edgecolor':'black'})

plt.title('Average Budget For Each Genres')
plt.tight_layout()
plt.show()
```
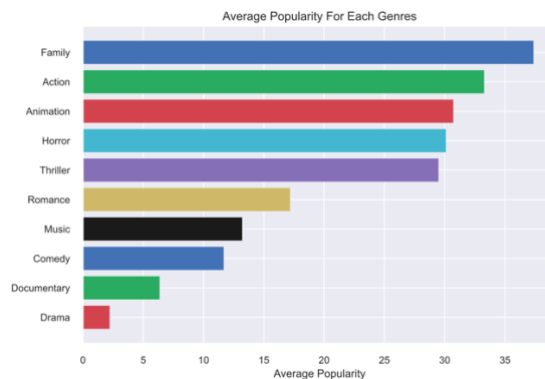
Average Budget For Each Genres



---

Results for the average of Profit, Budget and Popularity for each top 10 Genres:

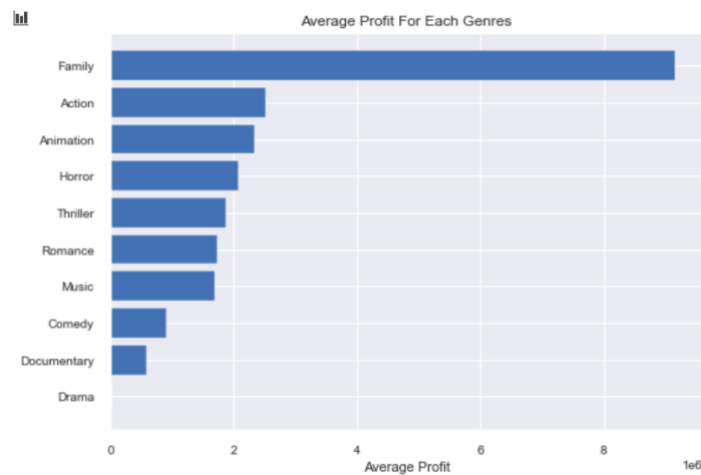Also, for the other columns we can apply the same plotting.

```
my_colors = 'rgbkymc'
genres = df_topGenres['genres']
avgBudget = df_topGenres['avgProductionBudget'].sort_values()
plt.barh(genres, avgBudget, color= my_colors)

plt.title("Average Production Budget For Each Genres")
plt.xlabel("Average Production Budget")
plt.tight_layout()
plt.show()
```

Average Profit For Each Genres

Family is the genre where the more investment occurs in English movies due to its well-known popularity among people of any age. As a result, it is the most profitable genre in the industry in comparison with Action and Animation that are also well-recognized lately.

Result for comparing the average of Budget, Domestic Gross and Worldwide Gross for these top ten Genres in a single plot:

In this stage we want to visualize all three average "investment for a specific Genre" and its Gross in both worldwide and Domestic area in one figure:

```python
import numpy as np

genres = df_topGenres['genres']

x_indexes=np.arange(len(genres))
width=0.25

avgWorldwideGross = df_topGenres['avgWorldwideGross']
plt.bar(x_indexes - width ,avgWorldwideGross, width=width, label= 'Wolrdwide Gross', color='#444444')

avgDomesticBudget = df_topGenres['avgDomesticBudget']
plt.bar(x_indexes ,avgDomesticBudget, width=width, label= 'Domestic Gross', color='green' )


avgProductionBudget = df_topGenres['avgProductionBudget']
plt.bar(x_indexes + width ,avgProductionBudget, width=width, label= 'Production Budget')
#
plt.xticks(ticks=x_indexes, labels= genres)
plt.xlabel('Genre')
plt.ylabel('Amount (USD)')
plt.title('Average Budget and Gross(Domestic & Worldwide) by their Genre')
plt.legend()
plt.tight_layout()
plt.show()
```
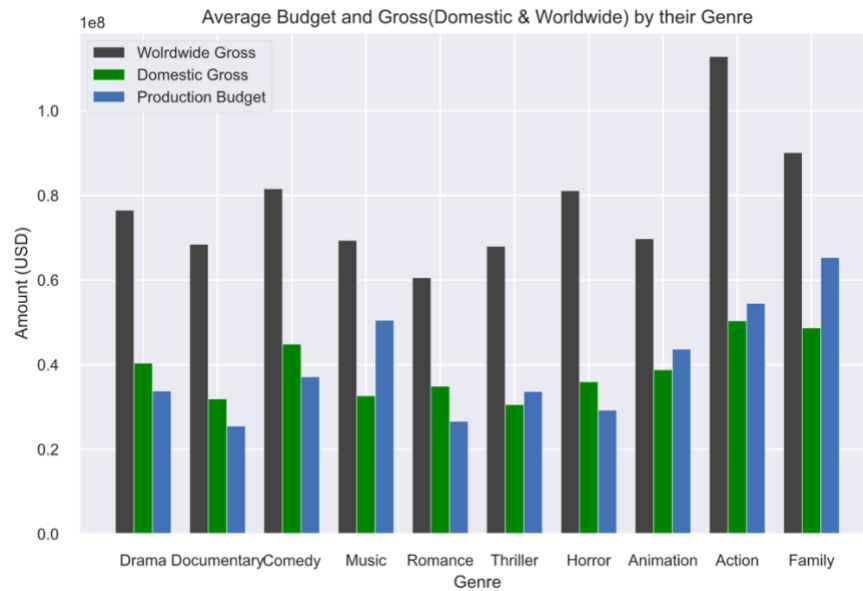
Average Budget and Gross(Domestic & Worldwide) by their Genre

In this plot we can see two main profitable genres: Action and Family. However, whereas the investment on the family Genre is more than Action, the later earns more Gross revenue than the former.

Third Part: (Visualization – Investigating on companies):

Loading Only Ten Top Company From Mongodb:

To have access to our mongodb and aggregate only ten top companies in a mongodb collection we need to write a function in a separate file as dbQueries.py as below:

```python
# Function to get top 10 production companies
def getTopProductionCompanies(topN=10):
    movies = list(collection.aggregate([{"$unwind": "$production_companies"}, {"$group": {"_id": "$production_companies",
                "avgRuntime": {"$avg": "$runtime"},
                "noMovies": {"$sum": 1}, "avgPopularity": {"$avg": "$popularity"},
                "avgVotesNo": {"$avg": "$vote_count"},
                "avgProductionBudget": {"$avg": {"$arrayElemAt": ['$aggregate.productionBudget', 0]}},
                "avgDomesticBudget": {"$avg": {"$arrayElemAt": ['$aggregate.domesticBudget', 0]}},
                "avgWorldwideGross": {"$avg": {"$arrayElemAt": ['$aggregate.worldwideGross', 0]}},
                "avgBudget": {"$avg": "$budget"},
                "avgRevenue": {"$avg": "$revenue"},
                "avgVotes": {"$avg": "$vote_average"}}},
                {"$sort": {"noMovies": -1, "avgRevenue": -1}}]))
    movies_json = []
    # In case we get less results than the specified number
    topN = (len(movies) if len(movies) < topN else topN)
    for i in range(topN):
        movie_json = {
            "production_companies": movies[i]["_id"],
            "noMovies": movies[i]["noMovies"],
            "avgPopularity": round(movies[i]["avgPopularity"], 3),
            "avgVotesNo": round(movies[i]["avgVotesNo"], 3),
            "avgBudget": round(movies[i]["avgBudget"], 3),
            # new fields from aggregation
            "avgProductionBudget": (0 if movies[i]["avgProductionBudget"] == None else round(movies[i]["avgProductionBudget"],3)),
            "avgDomesticBudget": (0 if movies[i]["avgDomesticBudget"] == None else round(movies[i]["avgDomesticBudget"],3)),
            "avgWorldwideGross": (0 if movies[i]["avgWorldwideGross"] == None else round(movies[i]["avgWorldwideGross"],3)),
            "avgRevenue": round(movies[i]["avgRevenue"], 3),
            "avgVotes": round(movies[i]["avgVotes"], 3)
        }
        movies_json.append(movie_json)
    return movies_json
```

Just like previous steps checks for any duplicates, null values or changing data types and then plot our

results. For avoiding any repetition, we did not write these data cleaning process here and just showing the results.

```
# Display results per production companies
# ======== Columns available: avgProductionBudget, avgDomesticBudget, avgWorldwideGross ========
re_json = dbQueries.getTopProductionCompanies(10)
df_topCompanies = pd.DataFrame.from_dict(re_json)
df_topCompanies.head(5)
```
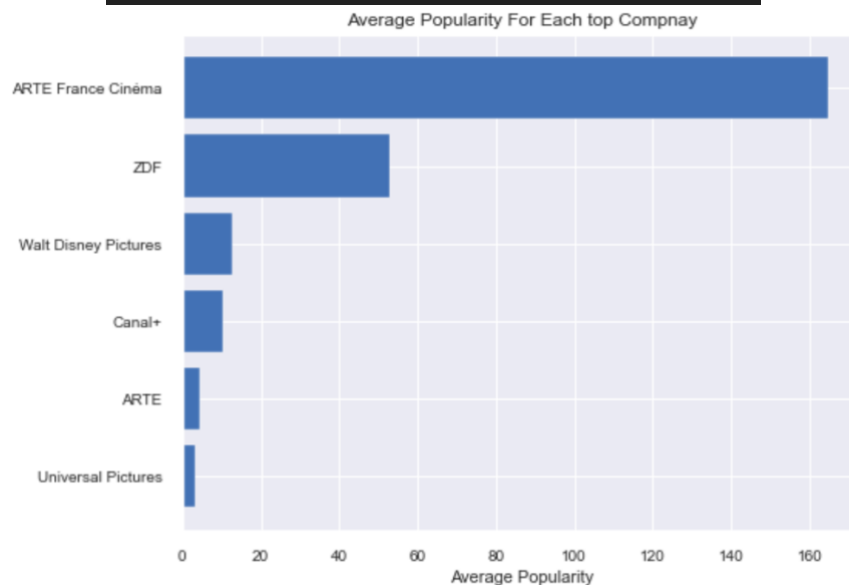
| | production_companies | noMovies | avgPopularity | avgVotesNo | avgBudget | avgProductionBudget | avgDomesticBudget | avgWorldwideGross | avgRevenue |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Universal Pictures | 25 | 52.762 | 1441.080 | 2.272000e+07 | 61400000.0 | 92396685.9 | 177364756.8 | 6.497059e+07 |
| 1 | BBC | 24 | 1.944 | 17.417 | 1.458333e+06 | 35000000.0 | 19389454.0 | 33422485.0 | 0.000000e+00 |
| 2 | ARTE | 22 | 3.155 | 70.273 | 8.540597e+05 | 12500000.0 | 4157491.0 | 45557491.0 | 1.820788e+06 |
| 3 | Canal+ | 18 | 10.390 | 183.278 | 2.009167e+06 | 12500000.0 | 4157491.0 | 45557491.0 | 2.024542e+07 |
| 4 | Walt Disney Pictures | 15 | 164.935 | 1476.267 | 5.826667e+07 | 94000000.0 | 41070064.8 | 75273052.8 | 3.746845e+07 |

## Results for the average of Popularity for each of these top companies:

Here we plot these top companies based on their order of average popularity and we do the same for the average production budget and profit and number of movies they have released, and we see if each company invest more and produce more movies will earn more profit and popularity.

```
plt.style.use('seaborn-notebook')
company = df_topCompanies['production_companies']
avgBudget = df_topCompanies['avgPopularity'].sort_values()
plt.barh(company, avgBudget)

#plt.xscale('log')
plt.title("Average Popularity For Each Compnay")
plt.xlabel("Average Popularity")
plt.tight_layout()
plt.show()
```

Results for comparing the average of Budget, Domestic Gross and Worldwide Gross for these top companies in a single plot:

```python
import numpy as np

companies = df_topCompanies['production_companies']

x_indexes=np.arange(len(companies))
width=0.25

avgWorldwideGross = df_topCompanies['avgWorldwideGross']
plt.bar(x_indexes - width ,avgWorldwideGross, width=width, label= 'Wolrdwide Gross', color='#444444')

avgDomesticBudget = df_topCompanies['avgDomesticBudget']
plt.bar(x_indexes ,avgDomesticBudget, width=width, label= 'Domestic Gross', color='green' )

avgProductionBudget = df_topCompanies['avgProductionBudget']
plt.bar(x_indexes + width ,avgProductionBudget, width=width, label= 'Production Budget')
#-----------------------------------------------#
plt.xticks(ticks=x_indexes, labels= companies)
plt.xlabel('Companies')
plt.ylabel('Amount (USD)')
plt.title('Average Budget and Gross(Domestic & Worldwide) by Top 10 Companies')
plt.legend()
plt.tight_layout()
plt.show()
```
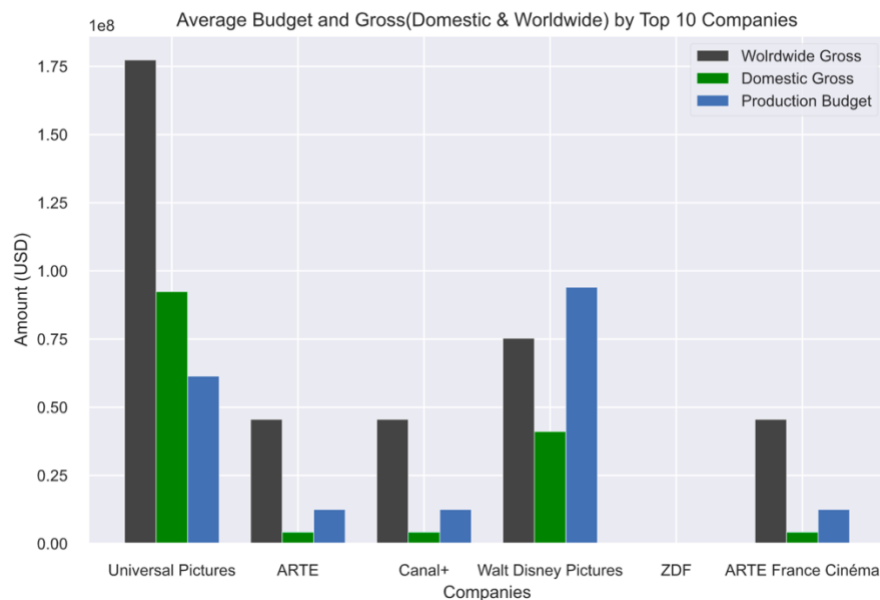


Here also we see that despite of Walt Disney invests more than Universal Pictures, Universal Pictures is still earning more gross revenue than its competitor.

## Conclusion:

In the end, by visualizing our data we can understand that one of the key features for the movie's success would be a higher investment in the budget to get more votes between people and as the result be more popular to achieve more revenue. We also see that key features for the popularity of the genres or success of the companies to have more profit could be the number of movies released or their higher average investment. During the project we got a good opportunity to play around data with the help of functionalities like Kafka and Mongo DB. Python gave us a lot of flexibility with libraries like pandas and matplotlib to better visualize our data. But we did find there is a scope for **future enhancements** as following:

- Aggregation with flask or Metabase for better visualization.
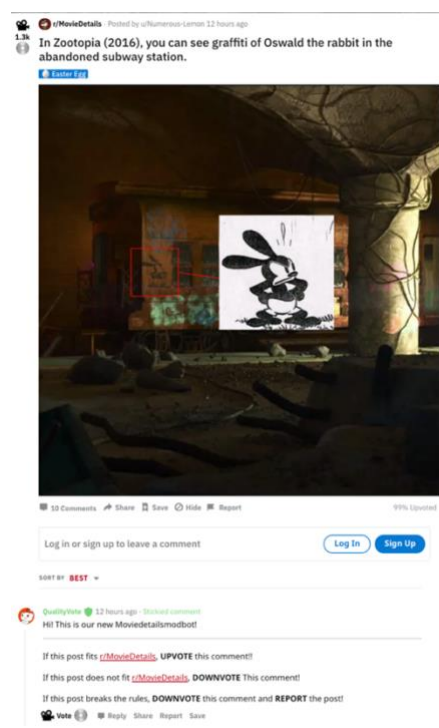- We can use docker instead of local machine to set this up.

## Additional Resources

**Requesting the Reddit API:**

The Reddit website do not provide developers by any API Key instead it gives auth, therefor for this purpose you also need to create an account in the reddit website. Then in the app link by entering your application name, an optional link and selecting the script for the personal use the secret key and the client id will be appear and ready for use.



In the reddit we have some subreddits, that everybody can join to his favorite one. For this project we consider working on one of these subreddits named as: MovieDetails to not only receive the title of the subreddit but also receive both comments and replies by their unique ids, we see here that the title of this subreddit and all its comments are fetching into our python program:

```python
import praw

reddit = praw.Reddit(
    client_id="yi_6eIYF8o2HLQ",
    client_secret="QgpuuStsvpvz91NGRBIjwl0KSjvL3g",
    password="reddit1234",
    user_agent="test",
    username="elnaz_di"
)

subreddit = reddit.subreddit('MovieDetails')
hot_python = subreddit.hot(limit=5)

# for submission in hot_python:
#     print(submission)

# for submission in hot_python:
#     print(submission.title)

for submission in hot_python:
    if not submission.stickied:
        print('Title : {}, ups : {}, downs : {}, Have we visited : {}'.format(
            submission.title, submission.ups, submission.downs, submission.visited))


# subreddit.unsubscribe()
# subreddit.subscribe()

comments = submission.comments.list()
for comment in comments:
    print(20*'-')
    print("Parent ID: ", comment.parent())
    print("Comment ID: ", comment.id)
    print(comment.body)
    if len(comment.replies) > 0:
        for reply in comment.replies:
            print("REPLY", reply.body)
```

```
/Users/dehkharghanielnaz/Desktop/movie/venv/bin/python "/Users/dehkharghanielnaz/Desktop/Project_Data Engineering/reddit.py"
Title : In Zootopia (2016), you can see graffiti of Oswald the rabbit in the abandoned subway station., ups : 1124, downs : 0,
Title : In The Last of the Mohicans (1992), as Chingachgook confronts Magua, Hawkeye stops Magua's remaining warriors from in
Title : In Final Destination (2000), there is a painting of a sword in Valerie Lewton's house, foreshadowing future events. E
Title : In Megamind (2010), you can see a picture of Mr Miyagi from the Karate Kid series., ups : 316, downs : 0, Have we vis
Title : In King Kong vs Godzilla (1962) the scene of Kong shoving a tree down Godzilla's throat is inspired by a publicity st
------------------------
Parent ID:  l5ks00
Comment ID:  gkursof
Hi! This is our new Moviedetailsmodbot!

----

If this post fits /r/MovieDetails, **UPVOTE** this comment!!

If this post does not fit /r/MovieDetails, **DOWNVOTE** This comment!

If this post breaks the rules, **DOWNVOTE** this comment and **REPORT** the post!
------------------------
```

## Docker Configuration- Apache Kafka using Confluent Platform:

### Step 1:

For using Apache Kafka through docker you need to have installed Docker then clone a confluent repository into your local system: https://github.com/confluentinc/cp-all-in-one Before using docker, you must be aware to set the memory from the default 2.00 GB to 8.00 GB to have a stable environment for working with Kafka Confluent. Then you can start Confluent by running a command in the terminal:
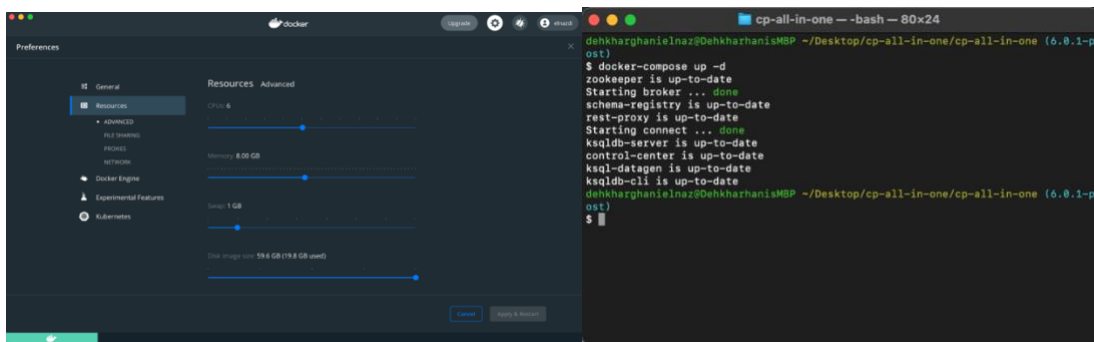
Figure 1.14 Starting Confluent in Docker. Left Image – set RAM to 8 GB. Right Image – starting Confluent.

here, you will see that Confluent is up and running in the Docker dashboard as well in the terminal:
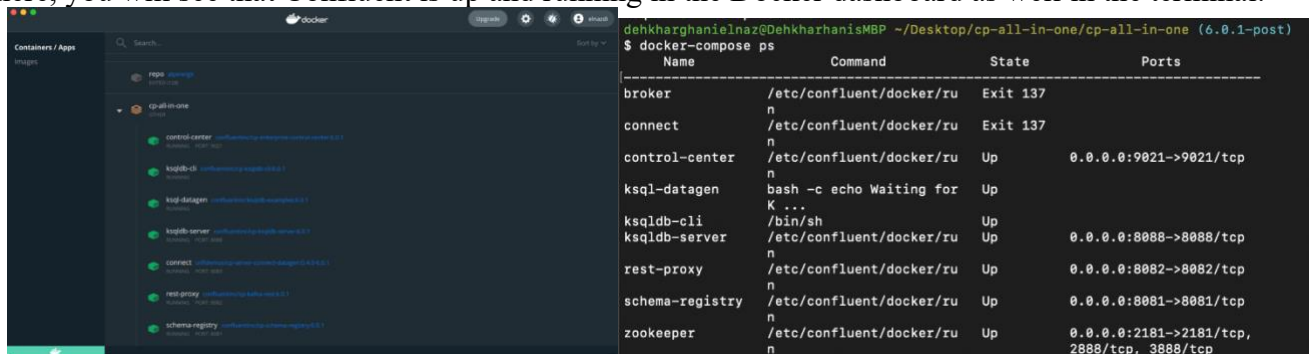


Figure 1.15 Confluent running. Left Image – Docker dashboard. Right Image – Terminal to see available services.

*Step 2:*

Create Kafka topics:

Now by using the Confluent Control Center we can create Kafka topics. Confluent Control Center provides the functionality for building and monitoring production data pipelines and event streaming applications. We Open the Control Center web interface at http://localhost:9021.
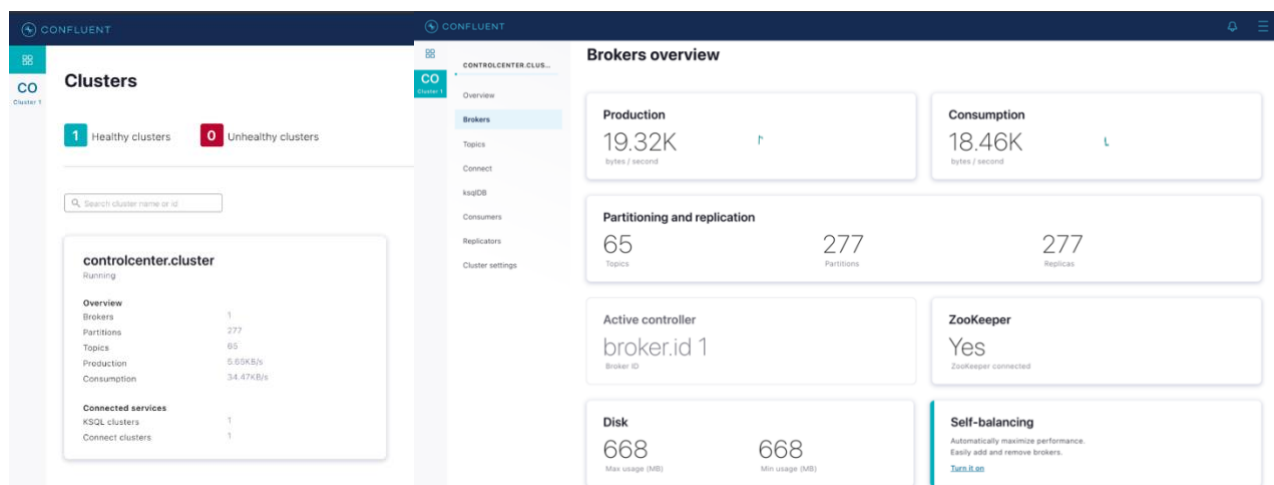


Figure 1.16 Confluent up and running. Left Image – cluster checkup. Right Image – Confluent dashboard

We continued our job with Confluent by adding topics, send and receiving message via python but since we could not stablish the MongoDB connector for this confluent we stopped here and continued our job with Kafka locally as an alternative solution.
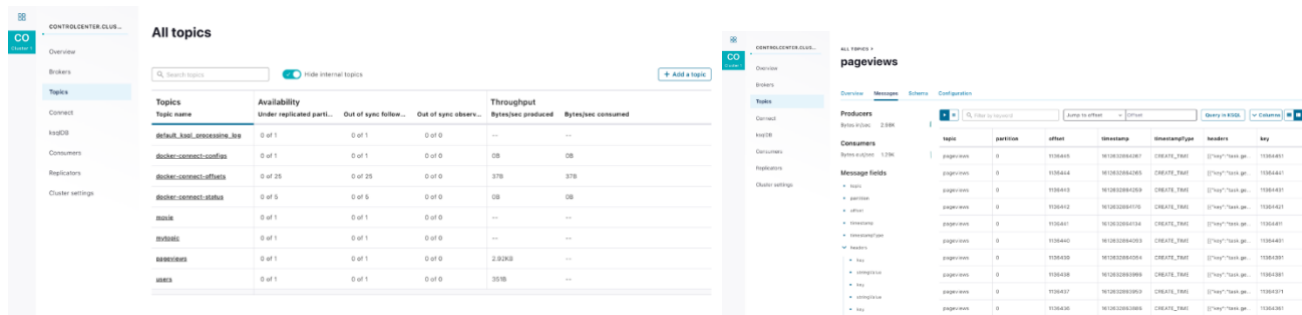
Figure 1.17 Topics in Confluent. Left Image – View all topics. Right Image – Create new topic.

## Resources:

[1] "The numbers", 1997-2021, Nash Information Service, LLC. https://www.the-numbers.com/movie/budgets/all

[2] "The Movie DB" (TMDB), 2008-2021, Community movie and TV database, https://www.themoviedb.org/

[3] Bat, A. "Investigating TMDB Movie Dataset", 2019, https://medium.com/my-data-camp-journey/investigating-tmdb-movie-datasets-4ee04c263915

[4] Waskon, M. "seaborn: statistical data visualization", 2012-2020, https://seaborn.pydata.org/index.html

[5] Richardson, L. "Beautiful Soup", (Last modification: October 3, 2020), https://www.crummy.com/software/BeautifulSoup/

[6] Cobrateam Revision, "Splinter", 2021, https://splinter.readthedocs.io/en/latest/

[7] Apache Sofwtare Foundation, "Apache Kafka", 2017, https://kafka.apache.org/quickstart

[8] Confluent Inc. "Quick Start for Apache Kafka using Confluent Platform (Docker)", 2021, https://docs.confluent.io/platform/current/quickstart/ce-docker-quickstart.html?utm_medium=sem&utm_source=google&utm_campaign=ch.sem_br.brand_tp.prs_tgt.confluent-brand_mt.mbm_rgn.emea_lng.eng_dv.all&utm_term=%2Bconfluent%20%2Bdocker&creative=&device=c&placement=&gclid=Cj0KCQiAmL-ABhDFARIsAKywVacHBIG8vBFffXpz12sVf0HtIgO7TqnrTfCQro5nwcMjHbhwqKvJPwEaAtZ7EALw_wcB

[9] MongoDB, "PyMongo 1.11.3 Documentation", 2018-2021, https://pymongo.readthedocs.io/en/stable/index.html