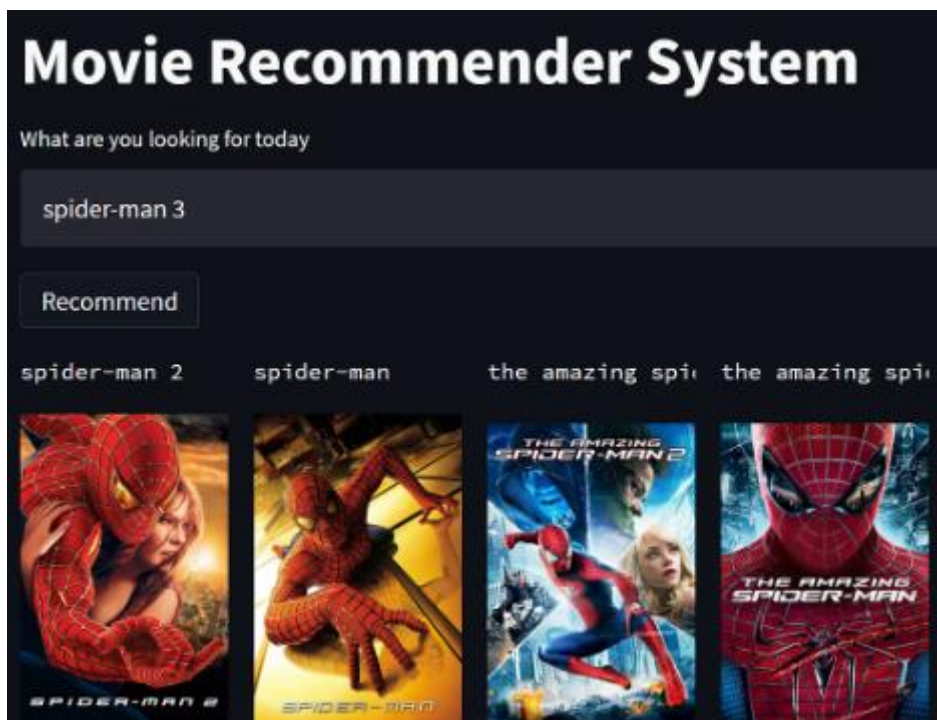# REPORT: MOVIE RECOMMENDATION SYSTEM
# Group Members:

Daniyar Abdrakhmanov, Doszhan Konysuly, Sunkar Kurmet

# 1. Introduction:
# Problem

The problem addressed in this work is the need for a recommendation system that suggests similar movies not only based on franchise but also on plot similarity. Traditional recommendation systems often rely on metadata such as genres, actors, or user ratings, which may not capture the nuanced similarities in movie plots. By leveraging natural language processing techniques and deep learning models, the aim is to provide more comprehensive and accurate movie recommendations that consider both franchise and plot elements.



# Literature Review:

Several approaches and solutions have been proposed in the field of recommendation systems and natural language processing:

1. **Content-Based Filtering:** Content-based filtering recommends items similar to those the user has liked in the past, based on item features. In the context of movies, this could involve analyzing movie summaries, keywords, or genres to identify similar films. The use of TF-IDF (Term Frequency-Inverse Document Frequency) and BERT (Bidirectional Encoder Representations from Transformers) models for feature extraction and similarity computation is common in content-based recommendation systems.

TF-IDF Vectorization

BERT Model

2. **Collaborative Filtering:** Collaborative filtering recommends items based on user behavior and preferences. This can include user ratings, viewing history, and social network information. While collaborative filtering is effective, it may not capture the semantic similarity between movies based on their plots.

3. **Hybrid Approaches:** Hybrid recommendation systems combine content-based and collaborative filtering techniques to provide more accurate and diverse recommendations. By integrating both approaches, hybrid systems aim to overcome the limitations of individual methods and enhance recommendation quality.

# Current Work:

The current work presents a novel approach to movie recommendation that combines content-based filtering with advanced natural language processing techniques. The workflow involves the following steps:

1. **Data Preprocessing:**

   - Movie data, including titles, overviews, and keywords, is retrieved from a CSV file.

   - Text data is preprocessed by converting to lowercase and removing punctuation.

2. **Feature Extraction:**

   - TF-IDF vectorization is applied to both movie overviews and keywords to represent textual information as numerical features.

- BERT model is utilized to extract semantic features from movie texts, capturing contextual information and relationships between words.

3. **Model Training:**

- Features extracted using TF-IDF and BERT are combined to create a comprehensive representation of each movie.

- Cosine similarity is calculated between movie features to measure similarity between movies based on both content and plot.

4. **Recommendation Generation:**

- Given a selected movie title, the system identifies similar movies based on cosine similarity scores.

- Recommendations are provided to the user, along with movie overviews, to facilitate informed decision-making.

DATA:

155000000,"[{"id": 10752, "name": "War"}, {"id": 36, "name": "History"}, {"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 18, "name": "Drama"}, {"id": 10749, "name": "Romance"}],,1966,"[{"id": 347, "name": "aristotle"}, {"id": 1160, "name": "egypt"}, {"id": 1200
150000000,"[{"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 10751, "name": "Family"}, {"id": 9648, "name": "Mystery"}],http://www.harrypotterorderofthephoenix.com/,675,"[{"id": 530, "name": "prophecy"}, {"id": 616, "name": "witch"}, {"id": 1014, "name": "loss
150000000,"[{"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 10751, "name": "Family"}],http://harrypotter.warnerbros.com/,674,"[{"id": 2343, "name": "magic"}, {"id": 3737, "name": "dying and death"}, {"id": 3872, "name": "broom"}, {"id": 3873, "name": "sorcere
150000000,"[{"id": 14, "name": "Fantasy"}, {"id": 28, "name": "Action"}],http://www.sonypictures.com/movies/hancock/,8960,"[{"id": 334, "name": "flying"}, {"id": 567, "name": "alcohol"}, {"id": 2036, "name": "love of one's life"}, {"id": 3691, "name": "forbidden love"}, {"id": 4663, "nan
150000000,"[{"id": 18, "name": "Drama"}, {"id": 27, "name": "Horror"}, {"id": 28, "name": "Action"}, {"id": 53, "name": "Thriller"}, {"id": 878, "name": "Science Fiction"}],http://iamlegend.warnerbros.com/,6479,"[{"id": 83, "name": "saving the world"}, {"id": 4190, "name": "lost civilisa
150000000,"[{"id": 12, "name": "Adventure"}, {"id": 35, "name": "Comedy"}, {"id": 10751, "name": "Family"}, {"id": 14, "name": "Fantasy"}],https://www.warnerbros.com/charlie-and-chocolate-factory,118,"[{"id": 212, "name": "london england"}, {"id": 494, "name": "father son relationsl
150000000,"[{"id": 16, "name": "Animation"}, {"id": 35, "name": "Comedy"}, {"id": 10751, "name": "Family"}, {"id": 14, "name": "Fantasy"}],http://disney.go.com/disneypictures/ratatouille/,2062,"[{"id": 90, "name": "paris"}, {"id": 380, "name": "brother brother relationship"}, {"id": 996
150000000,"[{"id": 28, "name": "Action"}, {"id": 80, "name": "Crime"}, {"id": 18, "name": "Drama"}],http://www2.warnerbros.com/batmanbegins/index.html,272,"[{"id": 486, "name": "himalaya"}, {"id": 779, "name": "martial arts"}, {"id": 849, "name": "dc comics"}, {"id": 853, "name"
150000000,"[{"id": 10751, "name": "Family"}, {"id": 16, "name": "Animation"}],http://www.madagascar-themovie.com,10527,"[{"id": 409, "name": "africa"}, {"id": 931, "name": "jealousy"}, {"id": 1691, "name": "dance"}, {"id": 1899, "name": "hunger"}, {"id": 2043, "name": "lion"}, {"i
150000000,"[{"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 28, "name": "Action"}, {"id": 35, "name": "Comedy"}, {"id": 10751, "name": "Family"}],http://www.nightatthemuseummovie.com,18360,"[{"id": 2598, "name": "museum"}, {"id": 5648, "name": "theodore
150000000,"[{"id": 12, "name": "Adventure"}, {"id": 28, "name": "Action"}, {"id": 53, "name": "Thriller"}, {"id": 878, "name": "Science Fiction"}],http://www.x-menorigins.com/,2080,"[{"id": 417, "name": "corruption"}, {"id": 1852, "name": "mutant"}, {"id": 2792, "name": "boxer"}, {"id
150000000,"[{"id": 12, "name": "Adventure"}, {"id": 28, "name": "Action"}, {"id": 53, "name": "Thriller"}, {"id": 878, "name": "Science Fiction"}],,605,"[{"id": 83, "name": "saving the world"}, {"id": 310, "name": "artificial intelligence"}, {"id": 312, "name": "man vs machine"}, {"id": 334
150000000,"[{"id": 16, "name": "Animation"}, {"id": 12, "name": "Adventure"}, {"id": 10751, "name": "Family"}],http://movies.disney.com/frozen,109445,"[{"id": 2011, "name": "queen"}, {"id": 4344, "name": "musical"}, {"id": 7376, "name": "princess"}, {"id": 10085, "name": "betrayal"
150000000,"[{"id": 12, "name": "Adventure"}, {"id": 28, "name": "Action"}, {"id": 53, "name": "Thriller"}, {"id": 878, "name": "Science Fiction"}],,604,"[{"id": 83, "name": "saving the world"}, {"id": 310, "name": "artificial intelligence"}, {"id": 312, "name": "man vs machine"}, {"id": 779
170000000,"[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}],http://marvel.com/thor,76338,"[{"id": 8828, "name": "marvel comic"}, {"id": 9715, "name": "superhero"}, {"id": 9717, "name": "based on comic book"}, {"id": 171783, "name": "ho
150000000,"[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 878, "name": "Science Fiction"}, {"id": 53, "name": "Thriller"}],http://www.madmaxmovie.com/,76341,"[{"id": 2964, "name": "future"}, {"id": 3713, "name": "chase"}, {"id": 4458, "name": "post-apocalypti
150000000,"[{"id": 53, "name": "Thriller"}, {"id": 9648, "name": "Mystery"}],http://www.angelsanddemons.com/,13448,"[{"id": 588, "name": "rome"}, {"id": 716, "name": "vatican"}, {"id": 818, "name": "based on novel"}, {"id": 1715, "name": "symbolism"}, {"id": 5950, "name": "christi
150000000,"[{"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 28, "name": "Action"}],http://thor.marvel.com/,10195,"[{"id": 1508, "name": "new mexico"}, {"id": 5149, "name": "banishment"}, {"id": 5539, "name": "shield"}, {"id": 8828, "name": "marvel comic"}, {"i
150000000,"[{"id": 16, "name": "Animation"}, {"id": 10751, "name": "Family"}, {"id": 12, "name": "Adventure"}, {"id": 35, "name": "Comedy"}],http://movies.disney.com/bolt,13053,"[{"id": 1939, "name": "hamster"}, {"id": 9963, "name": "kids and family"}, {"id": 18165, "name": "anima
150000000,"[{"id": 14, "name": "Fantasy"}, {"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 10751, "name": "Family"}, {"id": 35, "name": "Comedy"}],http://disney.go.com/disneypictures/gforce,19585,"[{"id": 156904, "name": "diy"}, {"id": 179431, "name": "duringc
150000000,"[{"id": 12, "name": "Adventure"}],http://www.wrathofthetitansmovie.org,57165,"[{"id": 1449, "name": "underworld"}, {"id": 2033, "name": "hades"}, {"id": 2035, "name": "mythology"}, {"id": 2036, "name": "greek mythology"}, {"id": 8985, "name": "zeus"}, {"id": 161170, "nan
150000000,"[{"id": 35, "name": "Comedy"}, {"id": 14, "name": "Fantasy"}],http://darkshadowsmovie.warnerbros.com,62213,"[{"id": 616, "name": "witch"}, {"id": 2883, "name": "imprisonment"}, {"id": 3133, "name": "vampire"}, {"id": 10541, "name": "curse"}, {"id": 11860, "name": "fi

CODE:

```python
data['overview'] = data['overview'].str.lower()
data['overview'] = data['overview'].str.replace('[^\w\s]','')

data['keywords'] = data['keywords'].str.lower()
data['keywords'] = data['keywords'].str.replace('[^\w\s]','')

tfidf_overview = TfidfVectorizer(stop_words='english')
tfidf_matrix_overview = tfidf_overview.fit_transform(data['overview'].fillna(''))

tfidf_keywords = TfidfVectorizer(stop_words='english')
tfidf_matrix_keywords = tfidf_keywords.fit_transform(data['keywords'].fillna(''))

tfidf_combined = hstack((tfidf_matrix_overview, tfidf_matrix_keywords))

tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
model = BertModel.from_pretrained('bert-base-uncased')

def extract_features(text):
    inputs = tokenizer(text, return_tensors="pt", max_length=512, truncation=True)
    outputs = model(**inputs)
    features = outputs.last_hidden_state.mean(dim=1).squeeze().detach().numpy()
    return features

movie_overviews = data['overview'].fillna('').tolist()
movie_keywords = data['keywords'].fillna('').tolist()
movie_texts = [overview + ' ' + keywords for overview, keywords in zip(movie_overviews, movie_keywords)]
movie_features = [extract_features(text) for text in movie_texts]

movie_features = torch.tensor(movie_features)
cosine_sim = linear_kernel(movie_features, movie_features)

def recommend_similar_movies(movie_title, num_recommendations=5):
    movie_index = data[data['original_title'] == movie_title].index[0]
    similar_movies = list(enumerate(cosine_sim[movie_index]))
    similar_movies = sorted(similar_movies, key=lambda x: x[1], reverse=True)[1:num_recommendations+1]
    similar_movie_indices = [i[0] for i in similar_movies]
    return data['original_title'].iloc[similar_movie_indices]
```

# 2. Data and Methods

# Information about the Data

We utilized a dataset from The Movie Database (TMDB), containing information about 4,807 movies.

Also conducted thorough analysis and preprocessing of the dataset before implementing the recommendation system. The analysis included:

```
[42] import pandas as pd
     import json
     import matplotlib.pyplot as plt

     from google.colab import drive
     drive.mount('/content/drive')

     file_path = '/content/drive/My Drive/DATASET/movies.csv'
     data = pd.read_csv(file_path)

     data = data.dropna(subset=['overview'])

     all_genres = []
     for genres_str in data['genres']:
         genres_list = json.loads(genres_str.replace("'", "\""))
         for genre in genres_list:
             all_genres.append(genre['name'])

     genres_df = pd.DataFrame(all_genres, columns=['genre'])

     plt.figure(figsize=(8, 8))
     genres_df['genre'].value_counts().plot(kind='pie', autopct='%1.1f%%', startangle=140)
     plt.title('Distribution of Movie Genres')
     plt.axis('equal')
     plt.show()

     overview_lengths = data['overview'].str.split().apply(len)
     plt.hist(overview_lengths, bins=20, color='skyblue')
     plt.xlabel('Overview Length')
     plt.ylabel('Frequency')
     plt.title('Distribution of Movie Overview Lengths')
     plt.show()
```
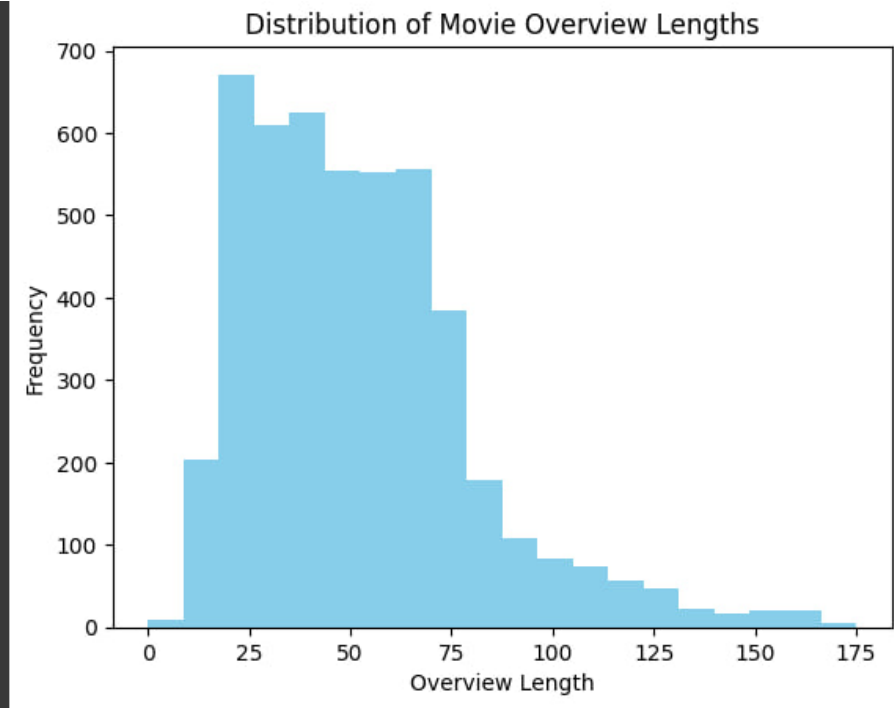
1. **Frequency and Length Analysis of Movie Overviews:**

   - We analyzed the distribution of overview lengths and frequencies to understand the descriptive patterns across movies.

   - A histogram was generated to visualize the distribution of overview lengths, providing insights into the variation in the length of movie descriptions.

Distribution of Movie Overview Lengths

2. **Genre Distribution Analysis:**

- We examined the distribution of genres present in the dataset to understand the diversity and prevalence of different movie genres.

- A bar chart was created to depict the count of each genre, enabling a clear visualization of the genre distribution in the dataset.



Distribution of Movie Genres

# Description of ML/DL Models Used

1. **TF-IDF Vectorization:**

   - We used TF-IDF (Term Frequency-Inverse Document Frequency) vectorization to represent textual data. Two TF-IDF vectorizers were employed, one for movie overviews and another for keywords.

2. **BertTokenizer and BertModel:**

   - We utilized the BERT (Bidirectional Encoder Representations from Transformers) model for extracting features from the combined textual data (overviews and keywords).

   - The **BertTokenizer** tokenizes input text, preparing it for model input.

   - The **BertModel** is used to extract contextualized word embeddings from the tokenized text.

3. **Linear Kernel for Cosine Similarity:**

   - We computed cosine similarity scores between the features extracted using BERT for each pair of movies.

**Theory Behind the Models**

- **TF-IDF Vectorization:**

  - TF-IDF is a numerical statistic that reflects the importance of a word in a document relative to a collection of documents.

  - It helps in capturing the significance of words in the corpus while filtering out common stopwords.

- **BERT (Bidirectional Encoder Representations from Transformers):**

  - BERT is a state-of-the-art natural language processing model developed by Google.

  - It utilizes a transformer architecture to capture bidirectional contextual relations in text data, providing rich word representations.

- **Cosine Similarity:**

  - Cosine similarity is a metric used to measure the similarity between two non-zero vectors in an inner product space.

  - In our context, it helps in quantifying the similarity between movie feature vectors extracted using BERT.

These models collectively allow us to recommend similar movies based not only on franchise but also on storyline similarities, enhancing the user experience in exploring related content.

# 3. Results:

Our endeavor to develop a movie recommendation system that suggests films based on both franchise and storyline similarities has yielded promising outcomes. The model effectively recommends movies that align with user preferences and interests, enriching the movie-watching experience.

Successes of the Model:

1. **Franchise and Storyline Similarities:**

   - The model successfully identifies movies not only based on franchise associations but also on thematic and storyline similarities.

2. **Enhanced User Experience:**

   - By considering both franchise and storyline elements, the recommendation system offers users a more personalized and engaging movie selection experience.

**Examples of Recommendations:**



These results demonstrate the efficacy and versatility of our recommendation system in providing tailored movie suggestions, thereby enhancing user satisfaction and engagement.

# 4. Discussion

# Critical Review of Results

While our recommendation system has shown promising results in suggesting movies based on franchise and storyline similarities, there are several aspects that warrant critical examination:

1. **Accuracy of Recommendations:**

   - Although the model provides recommendations, the accuracy of these suggestions may vary based on the complexity of movie plots and the diversity of user preferences.

   - Further refinement of the recommendation algorithm may be necessary to enhance the precision of movie suggestions.

2. **Handling Ambiguity and Subjectivity:**

   - The interpretation of movie plots and themes can be subjective, leading to varying perceptions of similarity among users.

   - Addressing ambiguity and subjectivity in movie recommendations requires a nuanced understanding of user preferences and context.

3. **Scalability and Performance:**

   - As the dataset expands and user demand grows, scalability and performance become crucial considerations.

   - Ensuring the efficiency and scalability of the recommendation system will be essential for maintaining optimal user experience under increased load.

4. **Feedback Mechanism:**

   - Incorporating a feedback mechanism that allows users to rate recommended movies can facilitate continuous improvement of the recommendation system.

   - Analyzing user feedback can provide valuable insights into the effectiveness of the model and areas for refinement.

# Next Steps

To address the aforementioned considerations and further enhance the effectiveness of our recommendation system, the following steps are proposed:

1. **Fine-tuning Model Parameters:**

- Conduct thorough experimentation and fine-tuning of model parameters to optimize the accuracy and relevance of movie recommendations.

- Explore techniques such as hyperparameter tuning and model ensemble methods to improve recommendation performance.

2. **Incorporating User Feedback:**

- Implement a feedback loop that allows users to provide ratings and feedback on recommended movies.

- Analyze user interactions and preferences to iteratively refine the recommendation algorithm and enhance user satisfaction.

3. **Exploring Advanced Techniques:**

- Investigate advanced machine learning and deep learning techniques, such as neural collaborative filtering and reinforcement learning, to capture complex user-item interactions and improve recommendation accuracy.

4. **Enhancing Data Quality and Diversity:**

- Continuously update and enrich the dataset with diverse movie titles and genres to reflect evolving user preferences and trends.

- Ensure data quality by conducting regular audits and validations to mitigate biases and inconsistencies.

5. **User-Centric Design:**

- Adopt a user-centric design approach to tailor recommendations based on individual user profiles, preferences, and viewing history.

- Provide users with customizable recommendation filters and preferences to personalize their movie discovery experience.

By embracing these next steps and addressing critical aspects of the recommendation system, we aim to deliver a more refined, accurate, and user-centric movie recommendation experience that resonates with a diverse audience.

# Sources

1. **TMDB Dataset:**

- The Movie Database (TMDB). Retrieved from https://www.themoviedb.org/

2. **Scikit-learn Documentation:**

- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine Learning in Python. Journal of Machine Learning Research, 12, 2825-2830. https://scikit-learn.org/stable/documentation.html

3. **Hugging Face Transformers Documentation:**

- Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., ... & Shleifer, S. (2019). Hugging Face's Transformers: State-of-the-art Natural Language Processing. ArXiv, abs/1910.03771. https://huggingface.co/transformers/

4. **PyTorch Documentation:**

- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., ... & Desmaison, A. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. Advances in Neural Information Processing Systems, 32. https://pytorch.org/docs/stable/index.html

5. **Pandas Documentation:**

- McKinney, W., & others. (2010). Data Structures for Statistical Computing in Python. In Proceedings of the 9th Python in Science Conference (Vol. 445, pp. 51-56). https://pandas.pydata.org/docs/

6. **Google Colab Documentation:**

- Google. (n.d.). Colaboratory. https://colab.research.google.com/notebooks/intro.ipynb