# Movie Recommender System: Project Briefing

## Executive Summary

This document provides a comprehensive overview of a movie recommender system designed to analyze film features and user preferences to suggest similar content. The project's primary objective is to enhance user engagement by helping them discover new movies aligned with their interests. The system employs a content-based/hybrid recommendation approach, utilizing a Hierarchical Model for clustering.Key findings from the data analysis indicate that **Director** and **Actor** are the most influential features for generating accurate recommendations. Prolific individuals like director Steven Spielberg and actor Samuel L. Jackson represent strong signals of user interest, and their stylistic signatures and career patterns provide rich data for the model. In contrast, **Cast Size** was found to be an unreliable feature for distinguishing between old and new films, though it can help differentiate between independent and blockbuster productions.A significant challenge identified is a major data quality issue: **over 20% of the dataset (1,023 movies) is missing producer information** . A smaller-scale issue exists with 30 movies lacking a director. The recommended strategy is to urgently enrich this data from external sources like IMDb or TMDB.Technically, the system is built on a robust MLOps pipeline using MLflow for experiment tracking, model packaging, validation, and deployment, ensuring a reproducible and production-ready workflow. The architecture is powered by a FastAPI framework for the API and a Streamlit interface for user interaction.

## 1. Project Objective and Recommendation Approach

The core goal of this project is to develop an intelligent movie recommendation system that suggests films based on a deep analysis of their features. The system focuses on artistic elements such as genre, director, writer, and cast to connect users with content that matches their viewing habits and preferences, thereby increasing platform engagement.The system utilizes a **content-based / hybrid recommendation approach** , leveraging a **Hierarchical Model** to group similar movies. This method focuses on the intrinsic properties of the movies themselves, making recommendations that are explainable and cluster-aware to improve both relevance and scalability.

## 2. Core Recommendation Features and Insights

Data analysis has identified several key creative roles as primary drivers for the recommendation engine. The system prioritizes features that strongly define a film's artistic identity and are closely followed by audiences.

### 2.1. Director as a Primary Feature

The director is identified as a powerful feature for predicting user preference due to their distinct stylistic signatures.

- **Insight:** The dataset contains 433 directors who have directed more than one movie. Steven Spielberg is the most prolific with 27 films, followed by Clint Eastwood (20), Martin Scorsese (20), and Ridley Scott (16).
- **Interpretation:** The presence of established and influential directors is a strong indicator of the dataset's quality. A director's style profoundly shapes a movie's identity (e.g., Tarantino with Crime/Action, Woody Allen with Comedy/Drama).

- **Recommendation:**
- Utilize **Director** as a primary feature in the recommendation model.
- Create **Director Signature Profiles** for those with 5+ movies, mapping their preferred genres, average ratings, and frequent collaborators.
- Implement **Director Style Clustering** to group filmmakers with similar aesthetics (e.g., Spielberg + Cameron, Nolan + Fincher).
- Recommend additional films from a director's catalog if a user shows affinity for 3+ of their movies.

## 2.2. Actor as a Primary Feature

Actors are a crucial feature, as users often follow specific performers. The system leverages actor data to model affinity and discover connections between films.

- **Insight:** Samuel L. Jackson is the most frequently appearing actor with 67 movies. An analysis of Leonardo DiCaprio's filmography shows a strong specialization in Drama (20 films, ~45% of his work), with versatility across Thriller, Romance, Crime, and Action, and a notable absence of Comedy.
- **Interpretation:** Popular actors appear in many films, and their career choices often establish a distinct "actor brand." Users may follow actors with specific genre specializations.
- **Recommendation:**
- Build an **Actor Similarity Matrix** based on shared genres and frequent co-appearances.
- Create **Actor-Genre Profiles** to recommend similar films based on an actor's dominant genres.
- Develop an **Actor Collaboration Network** to recommend films featuring frequent co-stars.
- Use a **Prolific Actors List** (30+ movies) as a "Popular Actor" signal to the model.

**Movie Counts for Prolific Actors**

| Actor | Movie Count |
| ------ | ------ |
| Samuel L. Jackson | 67 |
| Robert De Niro | 57 |
| Johnny Depp | 40 |
| Brad Pitt | 38 |
| Tom Hanks | 32 |
| Leonardo DiCaprio | 22 |

## 2.3. Crew Roles and Feature Prioritization

Analysis of the crew job distribution provides a clear hierarchy for feature importance.

- **Insight:** The most common crew role is Producer (10,206), followed by Executive Producer (6,177) and Director (5,166). Woody Allen is the most prolific writer with 24 screenplays, exemplifying the "Auteur" filmmaker who writes, directs, and acts.
- **Interpretation:** The high number of producers is due to multiple producer credits per film. The distribution of core creative and technical roles is balanced and realistic.
- **Recommendation:** Prioritize features based on their creative influence.
- **High Priority:** Director, Screenplay, Original Music Composer.
- **Medium Priority:** Producer, Editor, Director of Photography.
- **Low Priority:** Art Direction, Casting.
- Enrich crew data with awards and nominations to compute a "Technical Excellence Score".

While not as influential as primary creative roles, other data points provide valuable context for differentiating films.

*3.1 Cast Size Analysis*

- **Insight:** The average number of cast members per movie is 22.12. Critically, there is no meaningful difference in average cast size between old and new movies (a difference of only 0.37 actors).
- **Interpretation:** The film industry has maintained consistent production patterns regarding cast size over time. Technological advances like CGI have not significantly reduced the need for actors.
- **Recommendation:**
  - **Do not** use cast size to distinguish between old and new movies.
  - **Do** use cast size to differentiate movie types. The documentation suggests creating categories:
  - **Small (1–15 actors):** Often independent or drama-focused.
  - **Medium (16–30 actors):** Typical for most films.
  - **Large (31+ actors):** Often large-scale Action or Fantasy productions.
  - Focus on cast quality (actor popularity, awards) rather than quantity.

*3.2 Producer Analysis*

- **Insight:** Over 470 producers have produced more than one movie. The most prolific producers after the "Unknown" category are Tim Bevan (35 movies) and Scott Rudin (31 movies)).
- **Interpretation:** Major producers often specialize in certain genres or styles.
- **Recommendation:**
  - Use **Producer** as a secondary, less influential feature compared to Director.
  - Focus on producers with 10+ films to ensure a strong signal.
  - Perform **Production House Analysis** by grouping movies from major studios (e.g., Disney, Warner Bros., Universal) that have distinct production patterns.

## 4. Data Quality and Mitigation Strategies

The analysis revealed significant data integrity issues that must be addressed to ensure model accuracy.

- **Missing Producer Data (Major Issue):**
- **Insight:** 1,023 movies, representing over 20% of the entire dataset, have no listed producer ("None").
- **Interpretation:** This is a major data quality issue that can negatively impact recommendation accuracy.
- **Recommendation:** This issue requires urgent attention. The proposed solution is to enrich the dataset using external sources like IMDb or TMDB, or alternatively, to create a dedicated "Unknown Producer" category.
- **Missing Director Data:**
- **Insight:** 30 movies have no listed director.
- **Interpretation:** These are likely independent, documentary, or otherwise obscure films.

- **Recommendation:** Perform manual data enrichment where possible. Assign any remaining cases to an "Unknown Director" category. Consider excluding these movies from recommendations or lowering their priority in the system.

## 5. System Architecture, Performance, and Deployment

The project is supported by a modern technical stack and a rigorous MLOps methodology to ensure performance and reproducibility.

### 5.1. Model Performance and Results

- **Model:** A Hierarchical Model is used for clustering.
- **Feature Matrix:** The final shape of the feature matrix is (4803, 224).
- **Explained Variance:** The Explained Variance by Singular Value Decomposition (SVD) is **0.2557516136412949**.
- **Clustering Performance:**
- **Silhouette Score (2 clusters):** 0.5971
- **Silhouette Score (10 clusters):** 0.0836
- **Cluster Distribution:** The model identified two primary clusters, with one being overwhelmingly dominant:
- **Cluster 2:** 4661 movies
- **Cluster 1:** 142 movies
- **Example Recommendations:** Top similar movies identified by the model include *The Shining*, *1408*, and *The Beyond*.

### 5.2. MLOps Pipeline and Technical Stack

A full MLflow MLOps lifecycle has been implemented to manage the model from experimentation to production.

- **MLflow Implementation:** The pipeline includes:
- **MLflow Tracking:** For logging experiments and metrics.
- **MLflow Model Packaging:** For standardized model storage.
- **Model Signature:** To validate input/output schemas.
- **Quality Gate:** For automated model validation checks.
- **Transition to Production:** For deploying production-ready models.
- **MLflow Project:** For ensuring reproducible workflow execution.
- **Technical Stack:**
- **API:** FastAPI Framework
- **Interface:** Streamlit