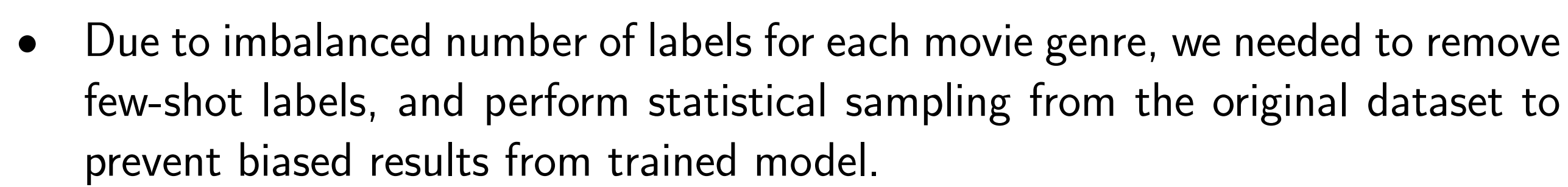


Amirhossein Soltani, Nigitha Selvarej

Project Overview:

Movie genre classification plays a crucial role in content recommendation systems by providing personalized movie suggestions, boosting user engagement. It significantly improves search functionality, allowing for efficient genre-based filtering and accurate search results. It also provides valuable insights into audience preferences and market trends. This information is crucial for content producers and marketers to understand which genres are gaining popularity, which helps in making informed decisions about future content creation and marketing strategies.

- 42,306 movie plot summaries extracted from IMDB and Wikipedia, including 362 different genres, metadata about the actors and the movies.



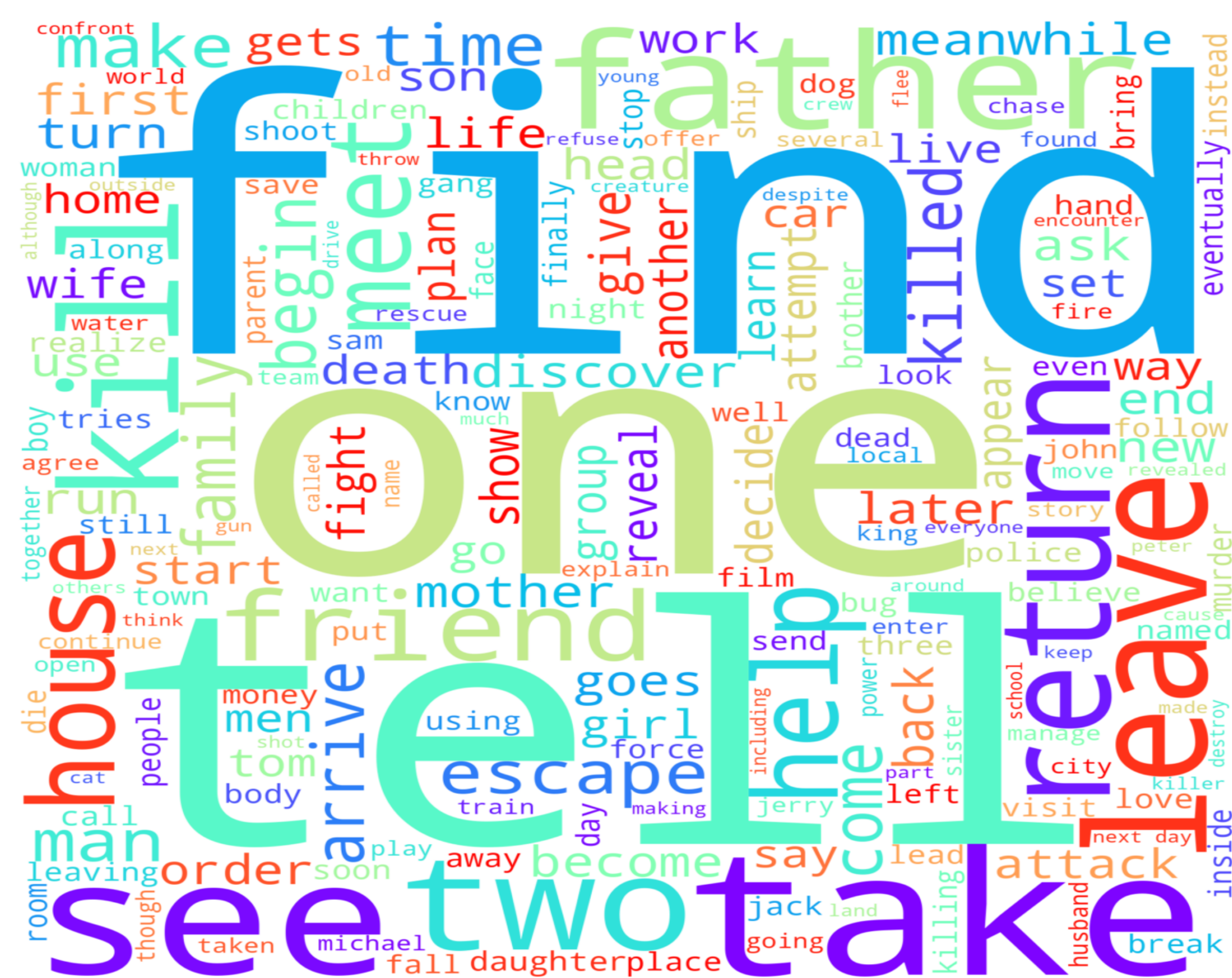
[1] Yinglong Ma, Xiaofeng Liu, Lijiao Zhao, Yue Liang, Peng Zhang, and Beihong Jin. Hybrid embedding-based text representation for hierarchical multi-label text classification. *Expert Systems with Applications*, 187:115905, 2022.

[2] Linkun Cai, Yu Song, Tao Liu, and Kunli Zhang. A hybrid bert model that incorporates label semantics via adjustive attention for multi-label text classification. *IEEE Access*, 8:152183–152192, 2020.

Your Name
Hochschule Bonn-Rhein-Sieg
Email: name.lastname@domain.com



In preprocessing step, we performed tokenization to convert text into tokens (words, sub-words, or characters) that the model can understand. Next steps included removing unnecessary characters, punctuation, or special symbols, normalization and lower casing, stemming and lemmatization. Next, we performed vectorization to encode text data into numerical format to capture the contextual features of the plot summaries and make them suitable for machine learning models. We utilized TF-IDF with ngram size (1,3) to capture more relationships between words. Finally, we trained our own language model as well as distinct classifiers and compared their results.



Your content here