

Amirhossein Soltani, Nigitha Selvarej

Project Overview:

Dataset:

-
- | Genre | Count |
|-----------------|-------|
| Drama | 34000 |
| Comedy | 16500 |
| Romance Film | 10500 |
| Black-and-white | 9000 |
| Action | 9000 |
| Thriller | 8800 |
| Short Film | 8200 |
| World cinema | 7200 |
| Crime Fiction | 7000 |
| Indie | 7000 |
| Documentary | 5800 |
| Horror | 5200 |
| Silent film | 5200 |
| Adventure | 5000 |
| Family Film | 4800 |

-
- | Genres | Counts |
|------------------|--------|
| Drama | 4200 |
| Comedy | 3400 |
| Romance Film | 2200 |
| Thriller | 2800 |
| Action | 3100 |
| Crime Fiction | 2400 |
| Horror | 3000 |
| Action/Adventure | 2800 |
| Adventure | 2800 |
| Family Film | 2600 |
| Short Film | 2900 |

Your content here

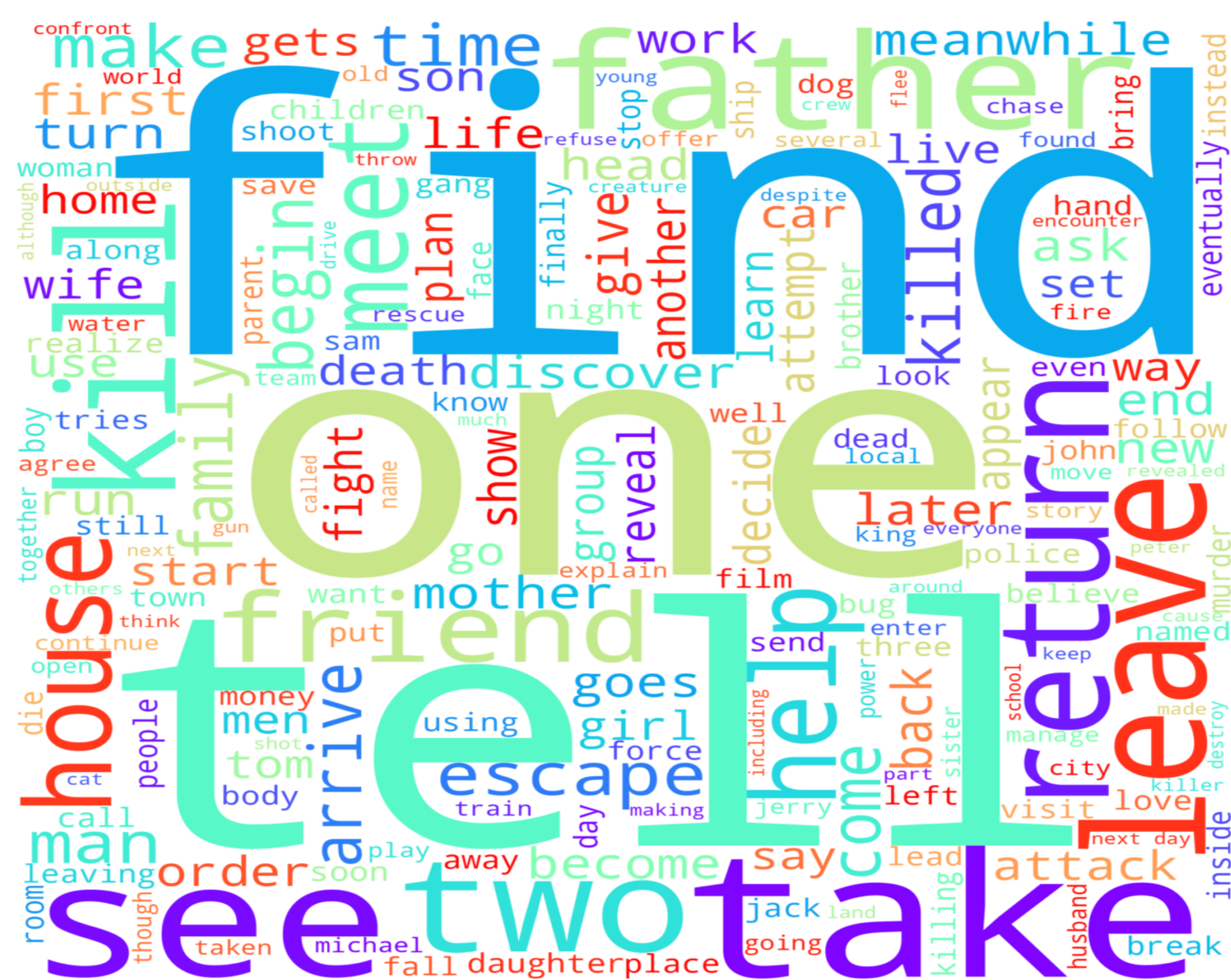
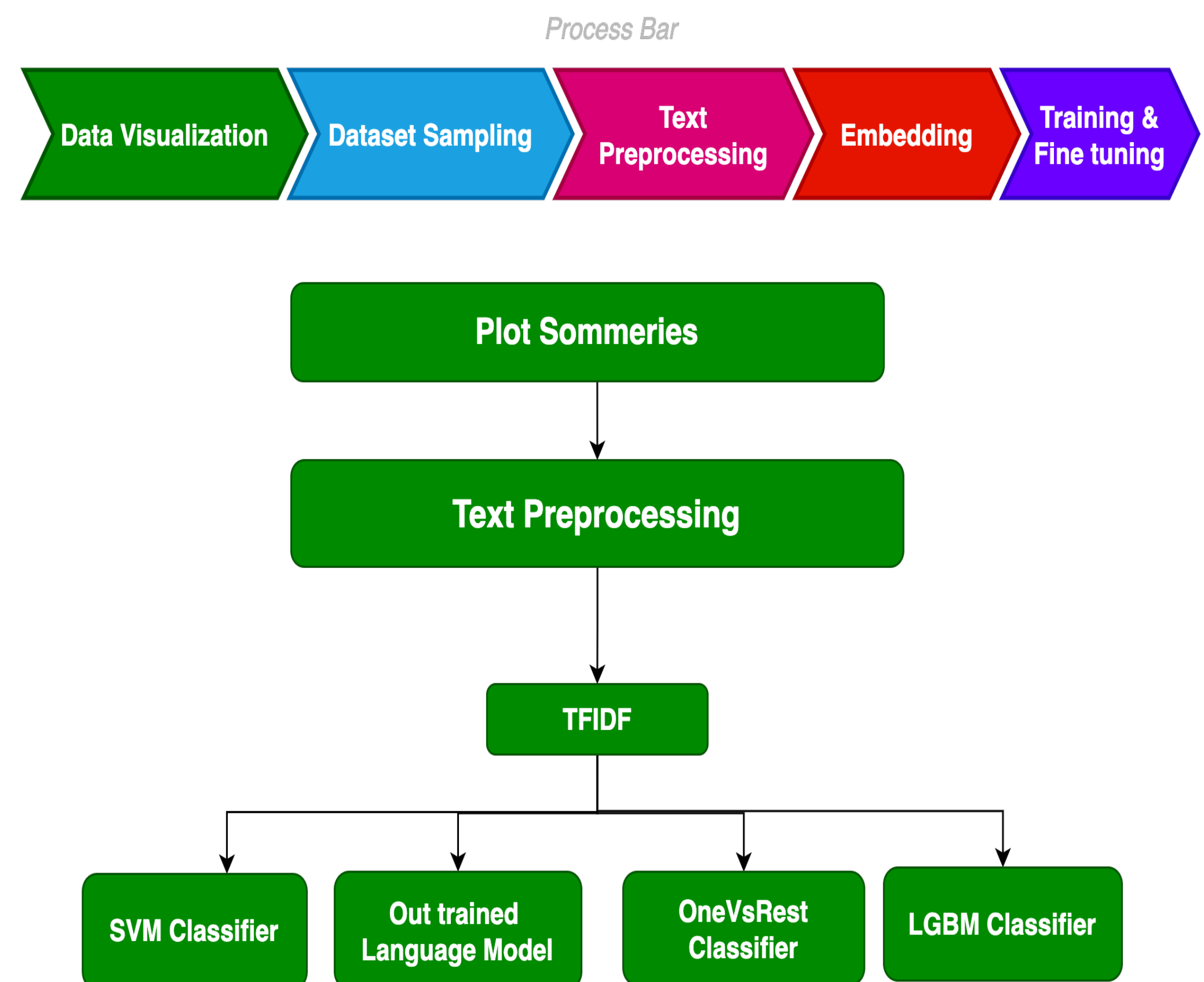
[1] Yinglong Ma, Xiaofeng Liu, Lijiao Zhao, Yue Liang, Peng Zhang, and Beihong Jin. Hybrid embedding-based text representation for hierarchical multi-label text classification. *Expert Systems with Applications*, 187:115905, 2022.

[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.

Nigitha Selvraj
Hochschule Bonn-Rhein-Sieg
Email: ??



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