

---

# Predicting Movie Genres Based on Plot Summaries

Team-13: AegiFi  
Akshit Sinha  
Adyansh Kakran

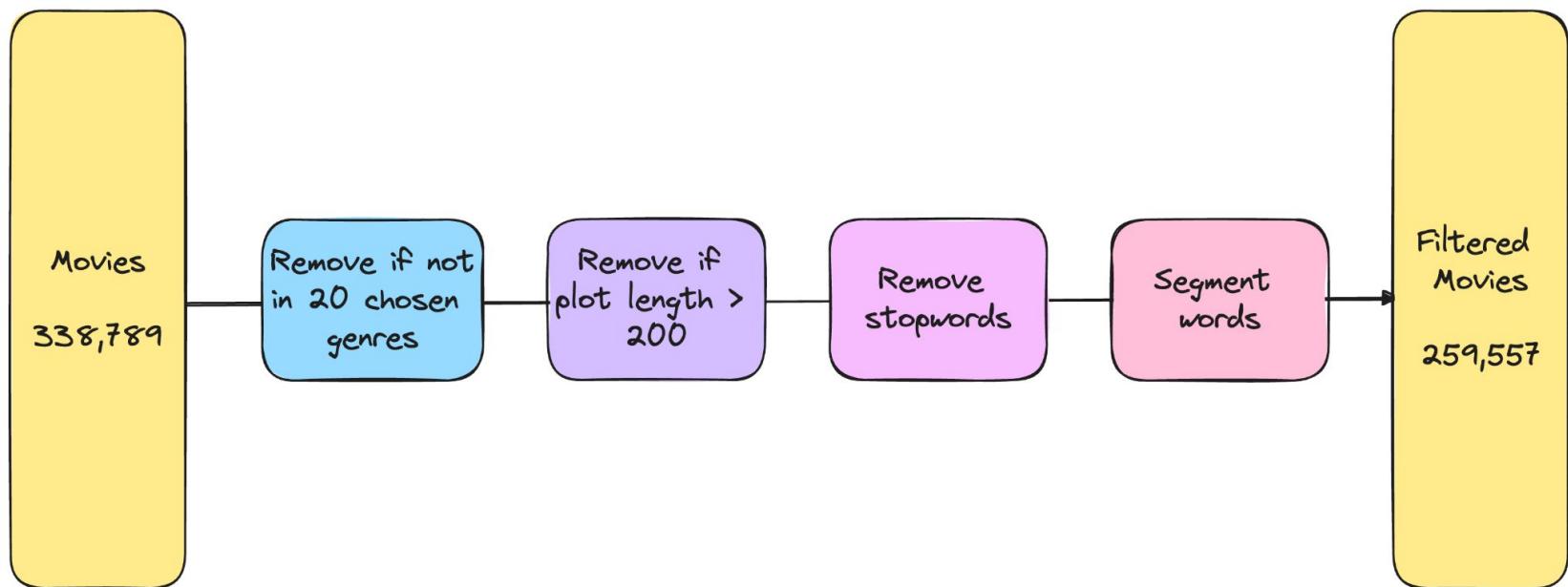


# Contents

1. Introduction to the paper and its implementation
2. Challenges faced while implementing the paper
3. Our work improving on the paper
4. Data and Model analysis

---

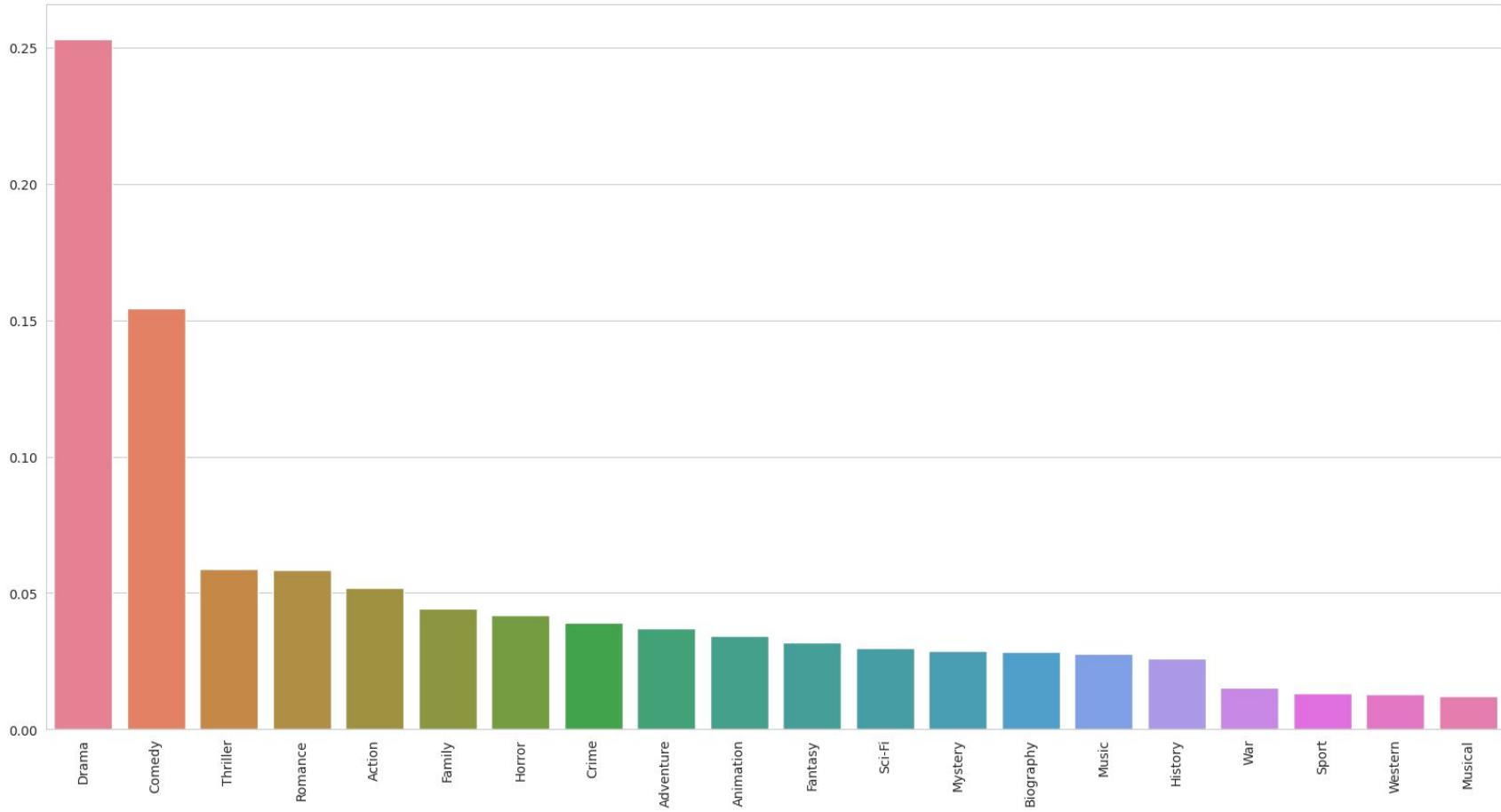
# Data preprocessing



---

# Exploratory Data Analysis

Distribution of Genres



# Word Cloud

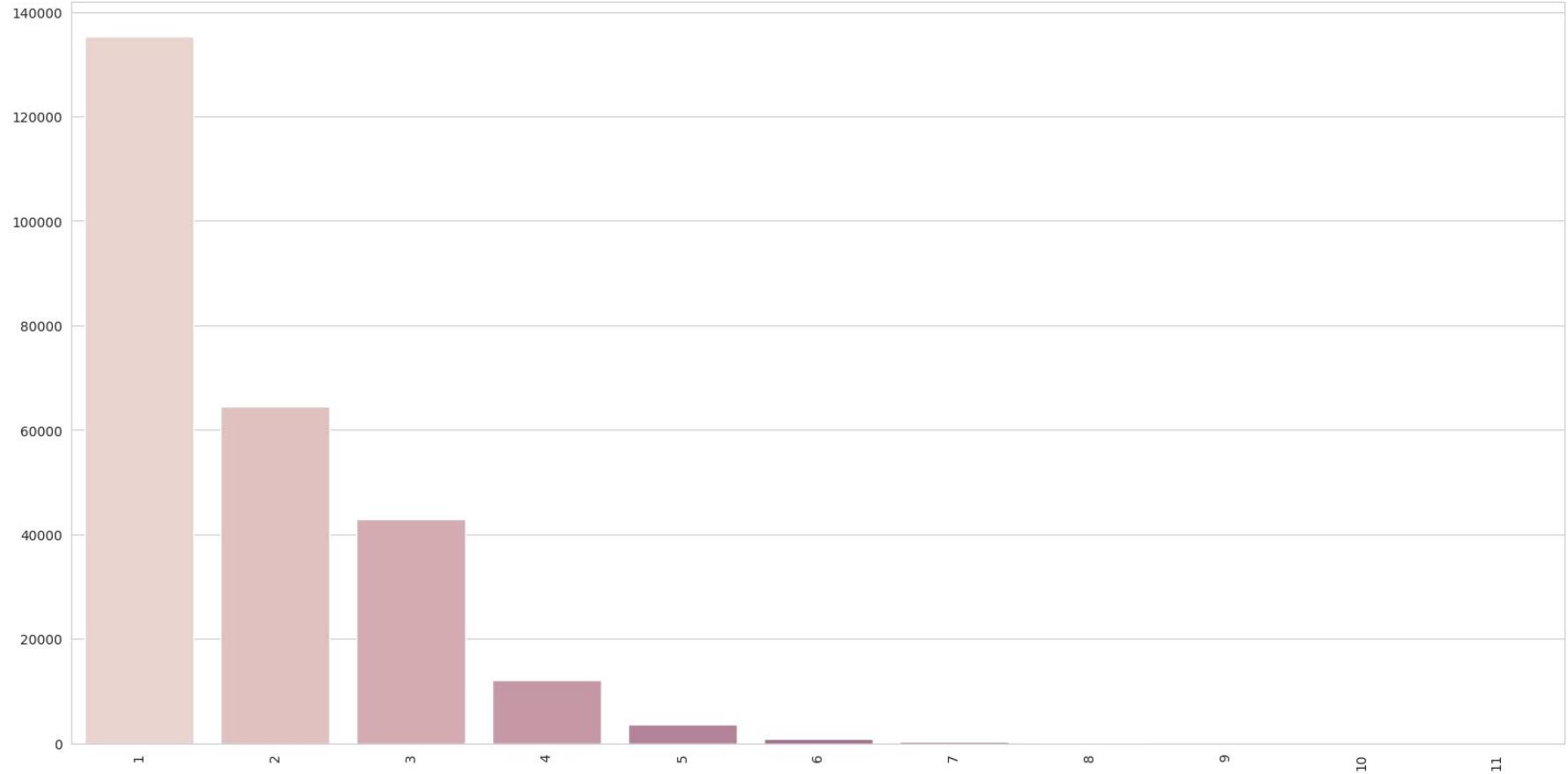
## Observations:

## Common words usually generic

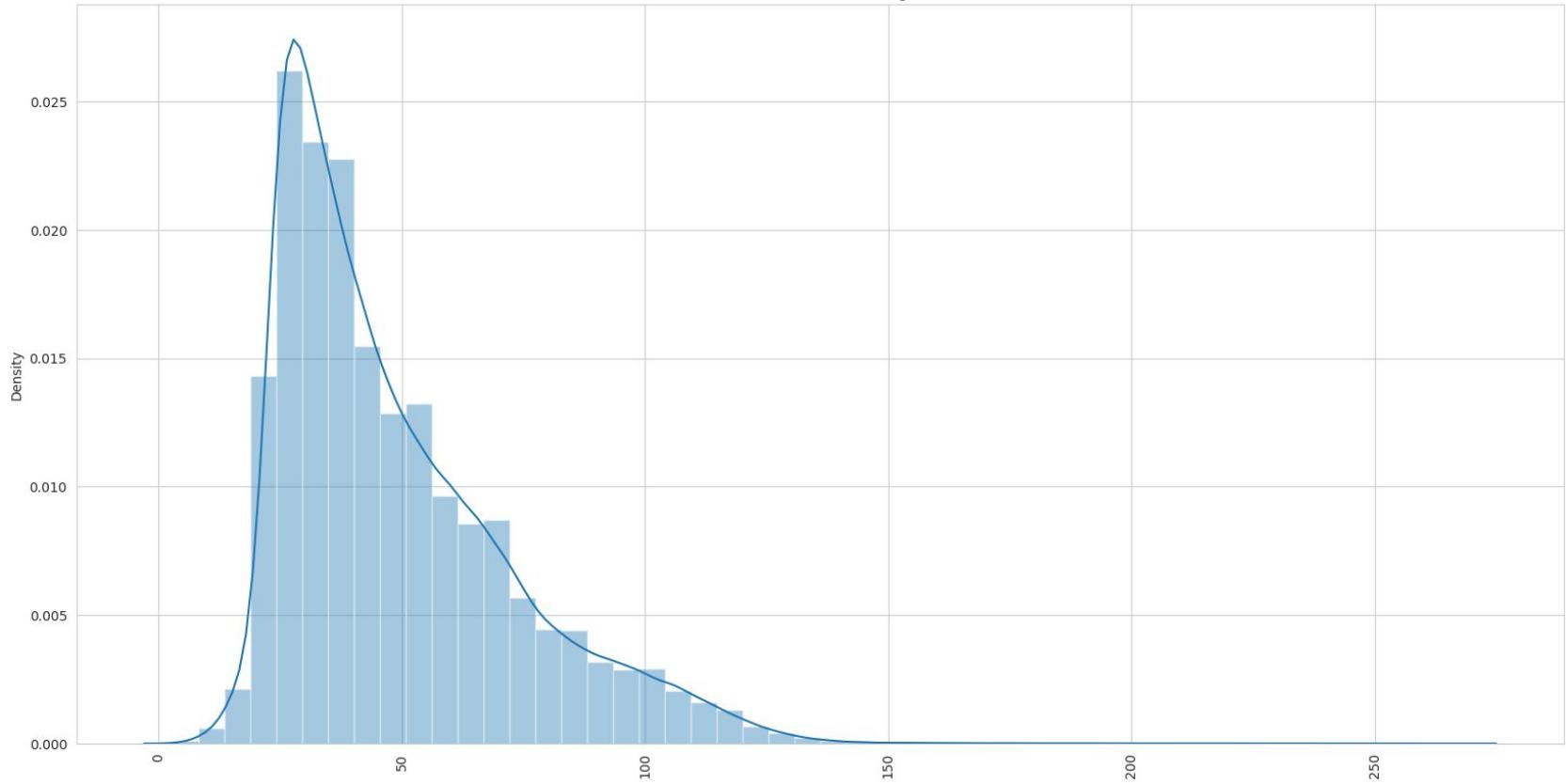
Some genre specific words:  
love, war, couple, etc



Distribution of Number of Genres



Distribution of Plot Length



---

# Machine Learning Models



# Naive Bayes

Two approaches used:

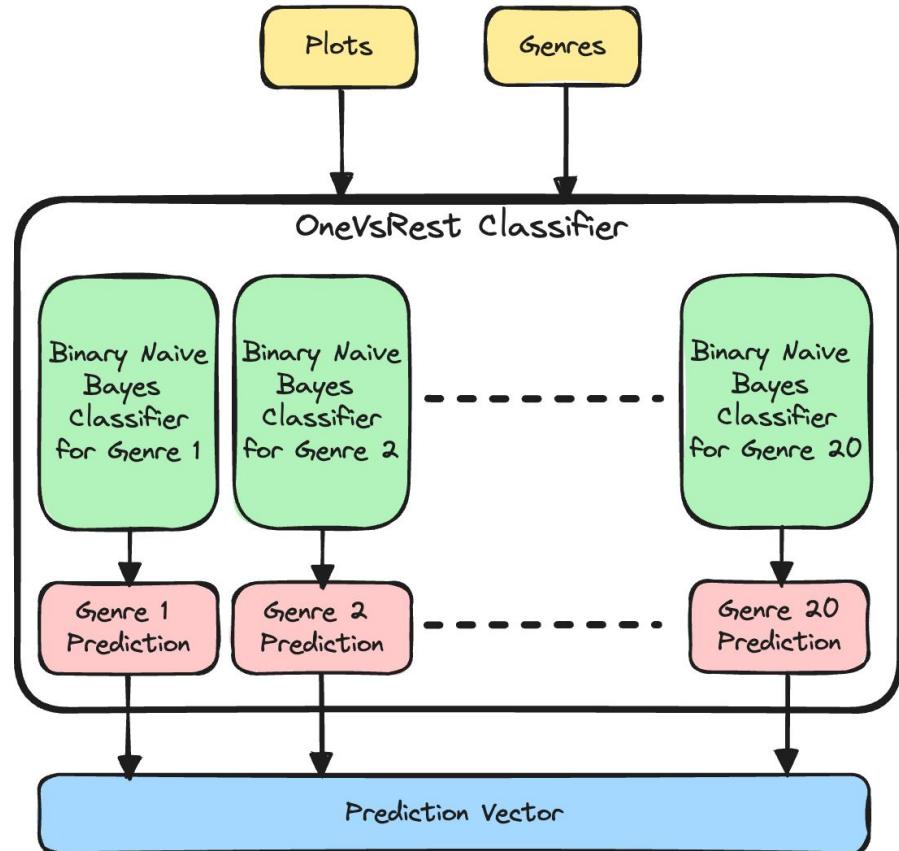
Bernoulli Naive Bayes using OneVsRest method

Construct 20 individual Binary Naive Bayes classifiers for each genre.

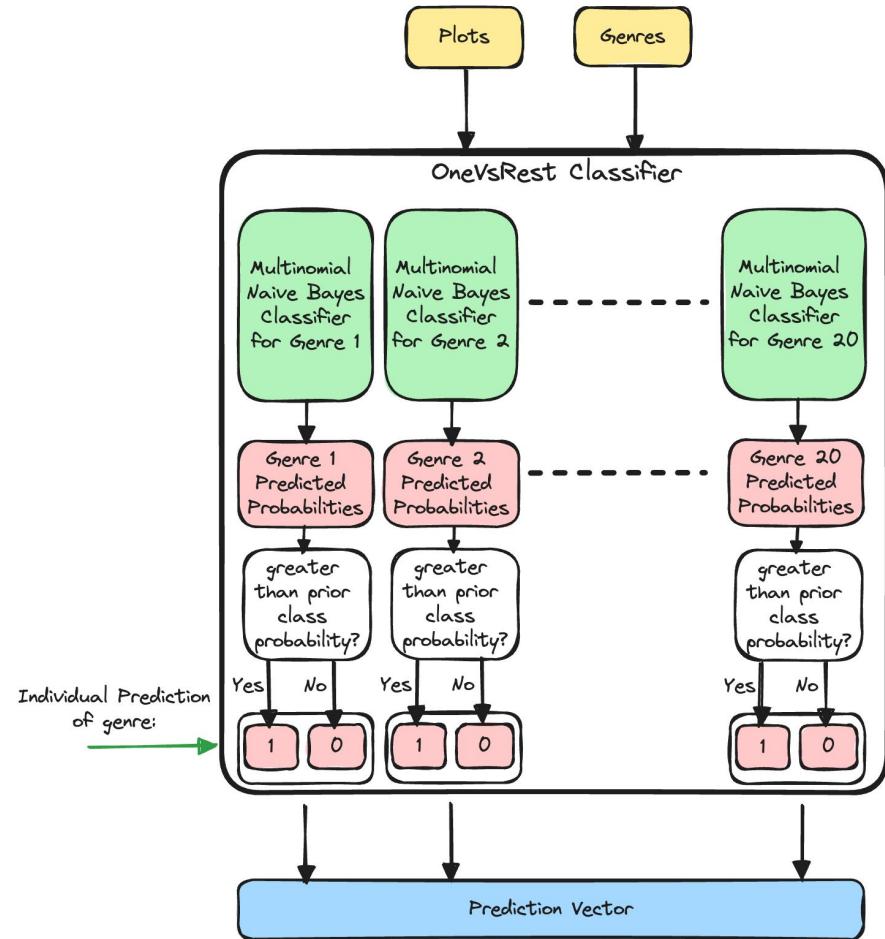
Multinomial Naive Bayes

Predict a genre if posterior probability is greater than prior probability

# Binary Naive Bayes



# Multinomial Naive Bayes





## XGBoost

Get pre trained word embeddings using the Google News Word2Vec corpus

Train an XGBClassifier on these embeddings. Get the output probabilities for the genres

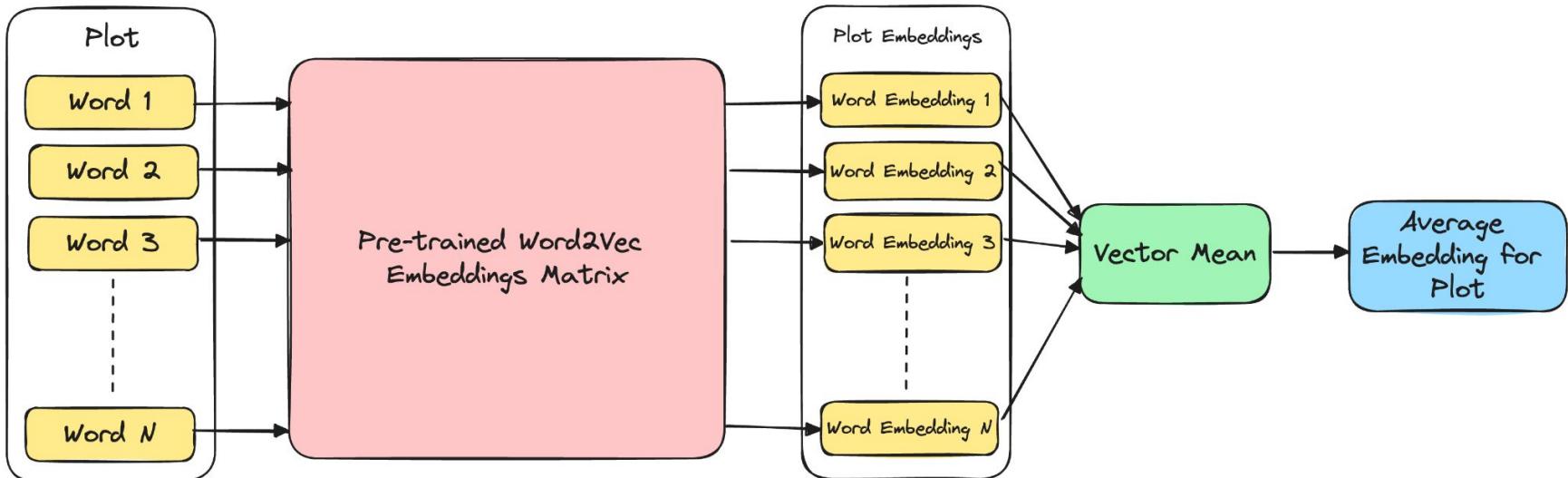
Calculate the threshold probability using the relation:

$$t(\mathbf{x}) = \arg \min_t | \{c \in C\} \text{ s.t. } P(c|\mathbf{x}) \leq t| + | \{c \in (\mathcal{Y} \setminus C)\} \text{ s.t. } P(c|\mathbf{x}) \geq t|$$

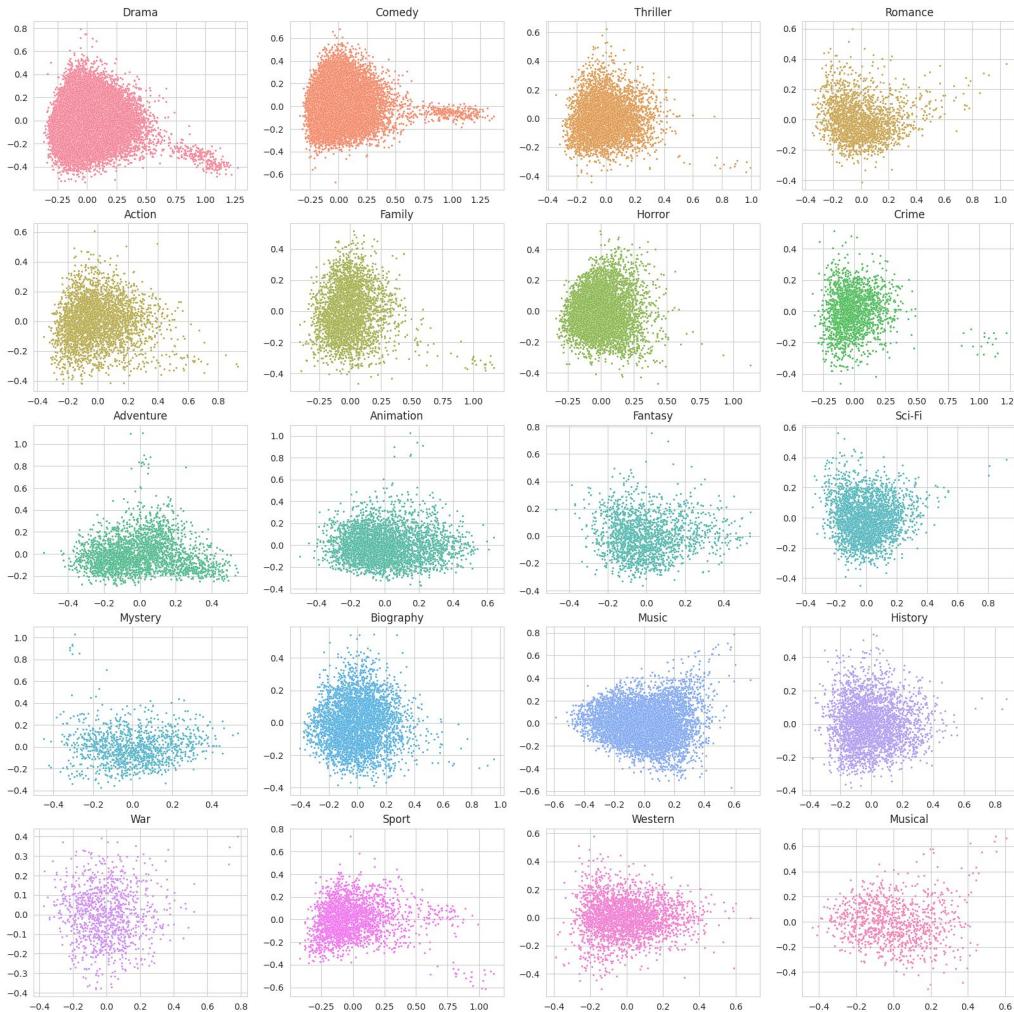
Feed the output probabilities as input and the calculated threshold as the label to an XGBRegressor to learn the best threshold values for the probabilities

---

# Word2Vec Embeddings

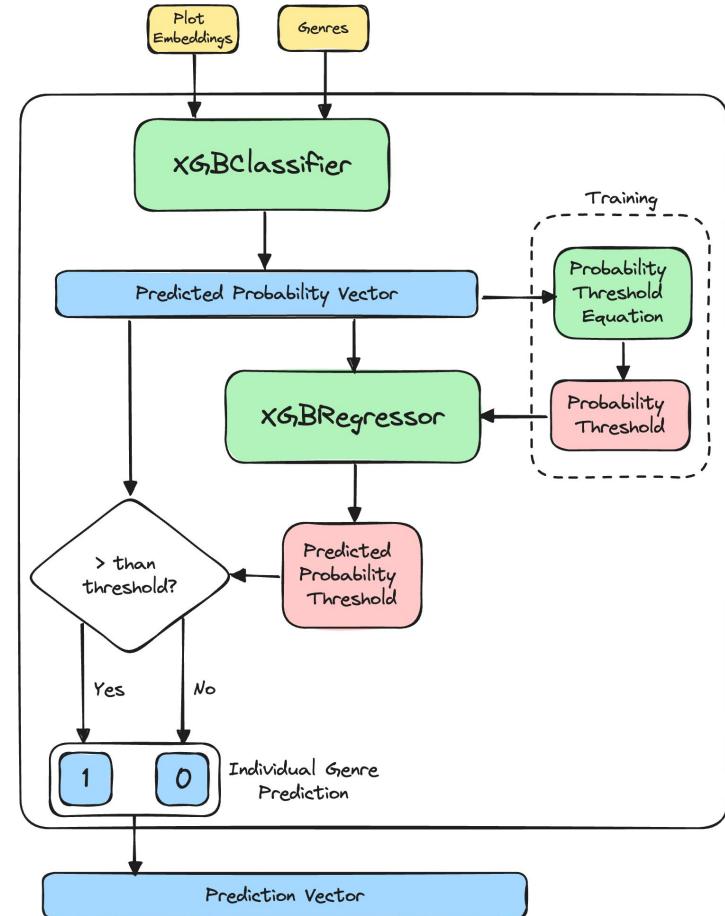


## Word2Vec Embeddings



# XGBoost Model

$$t(\mathbf{x}) = \arg \min_t |\{c \in C\} \text{ s.t. } P(c|\mathbf{x}) \leq t| + |\{c \in (\mathcal{Y} \setminus C)\} \text{ s.t. } P(c|\mathbf{x}) \geq t|$$





# GRU

3 approaches used

Binary GRU:

Output linear layer is fed to a sigmoid activation function, and genres with a sigmoid value  $> 0.5$  are predicted

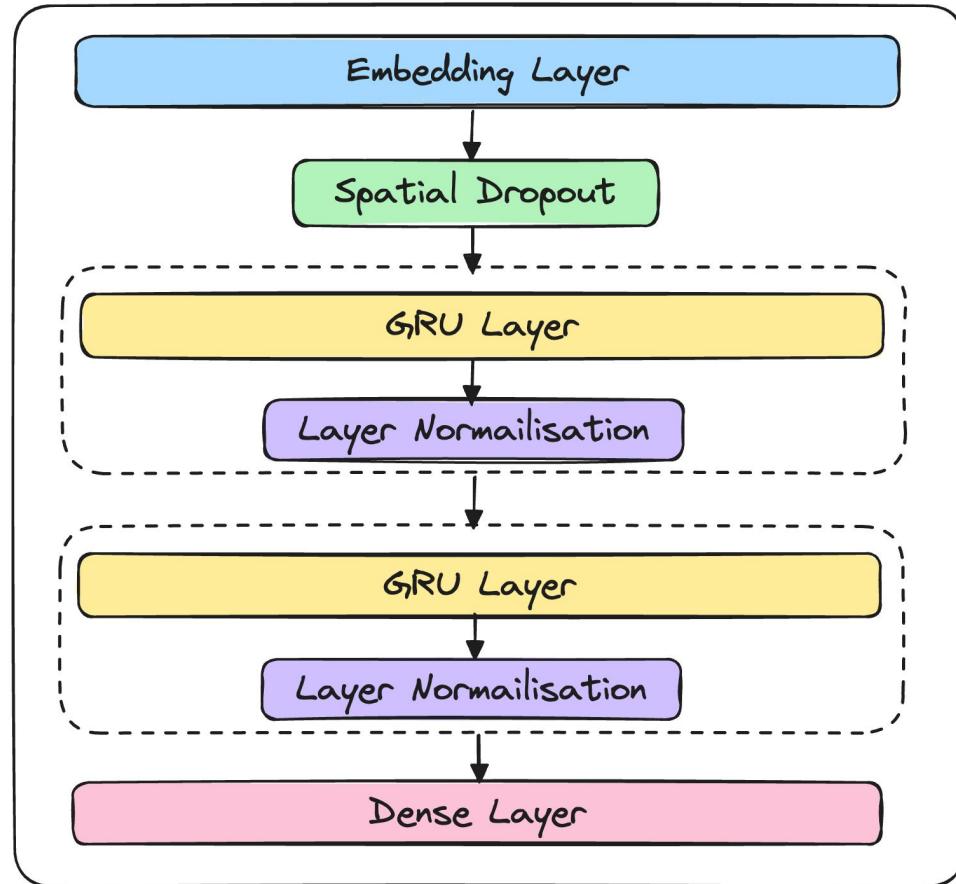
Rank GRU:

Linear layer outputs are interpreted as ranks of the genres. A sigmoid function is applied to convert them to values between 0 and 1, and best probability threshold is calculated as previously, and an XGBRegressor is trained on these values

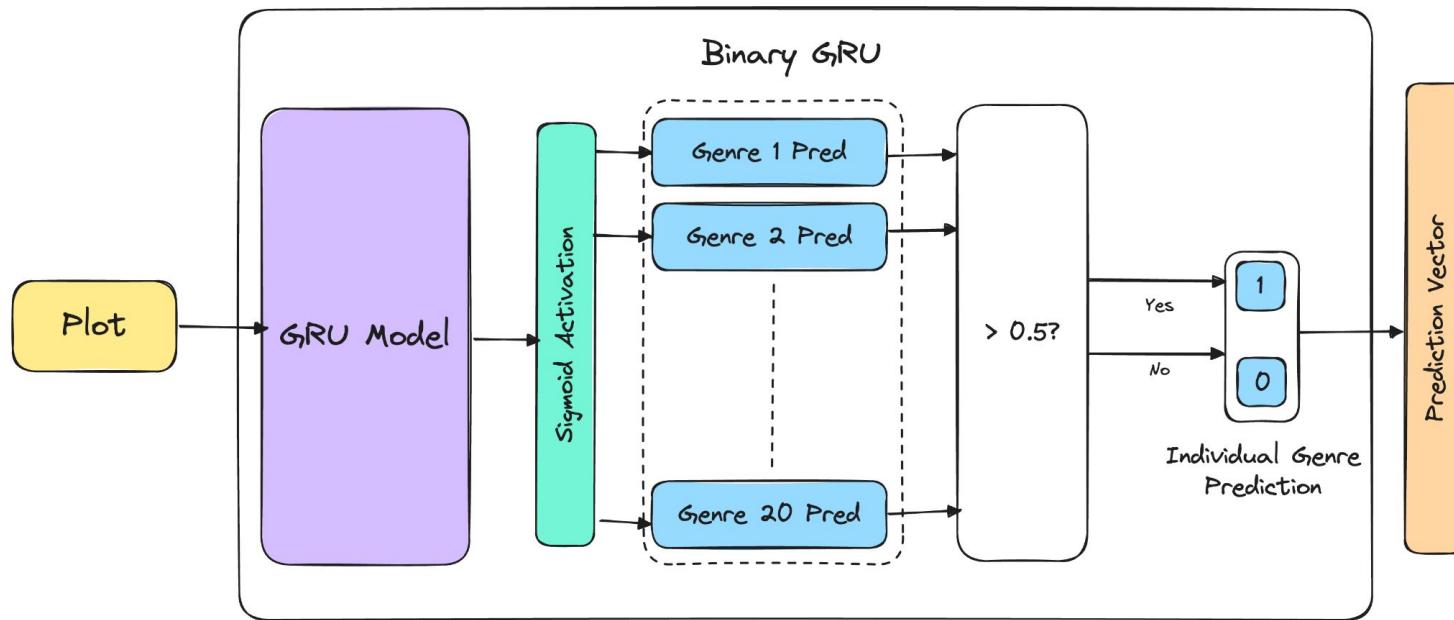
Multinomial GRU:

Outputs are fed into a softmax function to estimate probability of each genre. Further, the target probabilities of each genre are set to  $1/k$  if  $C$  is in the set of  $k$  correct genres, and 0 otherwise. Network is trained to minimise cross entropy between predicted and target categorical distribution. Again, threshold values are estimated using the previous method.

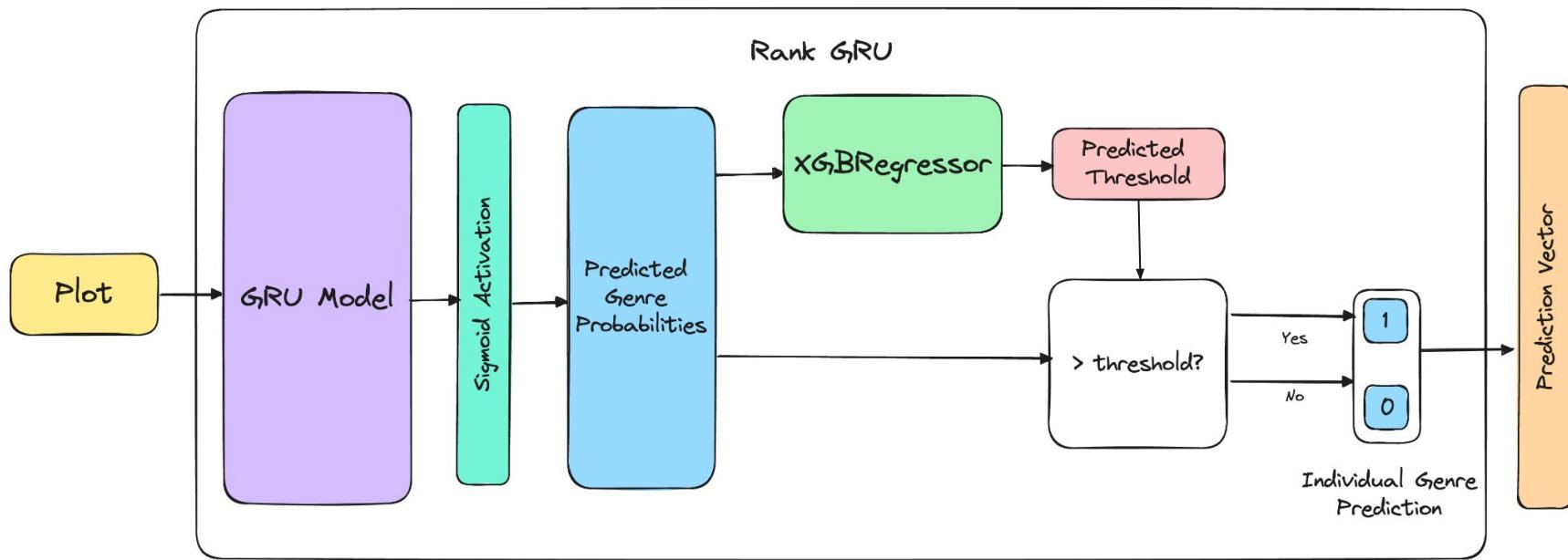
# GRU Architecture



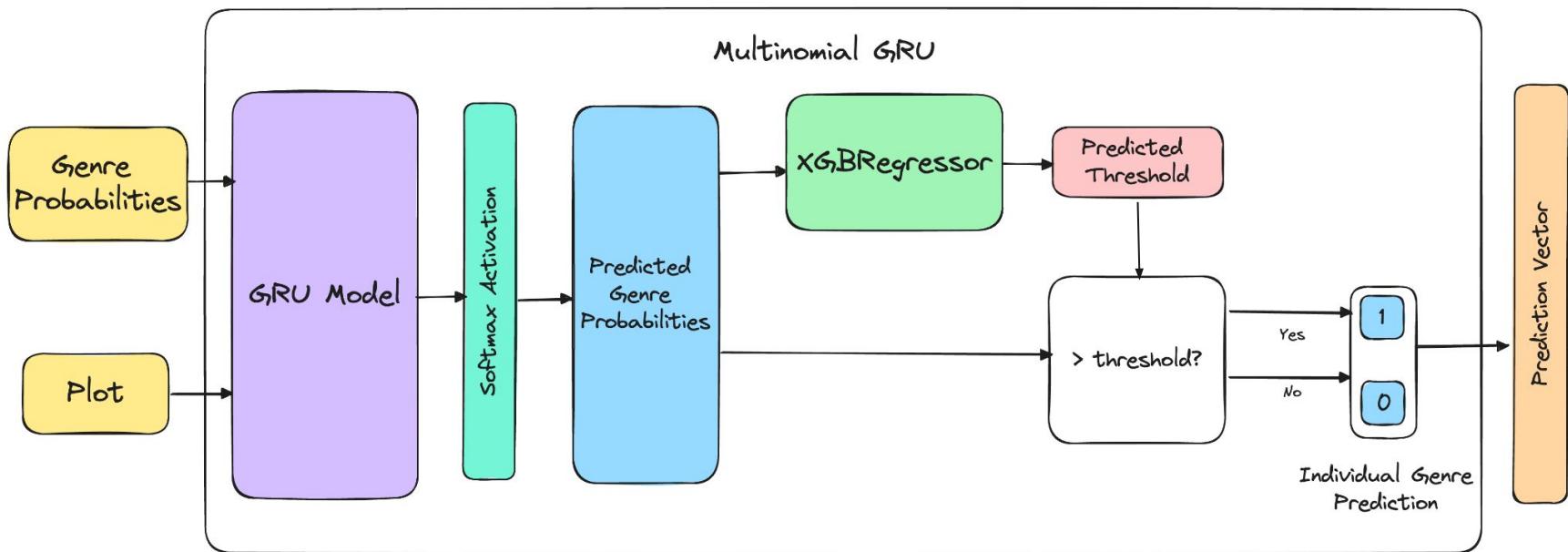
# Binary GRU



# Rank GRU



# Multinomial GRU



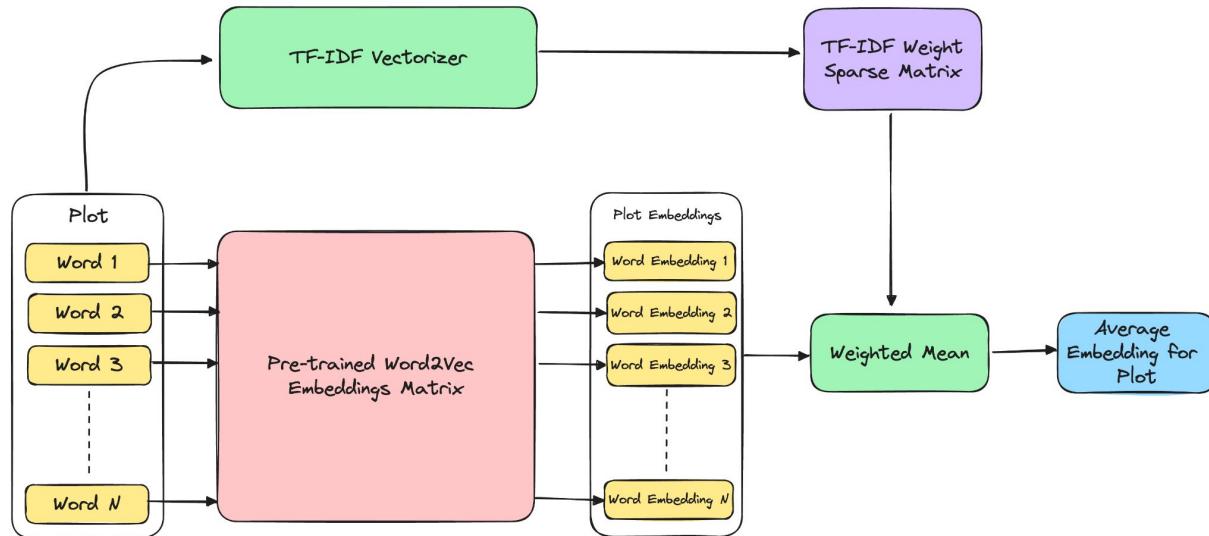
---

# Improvements over Publication

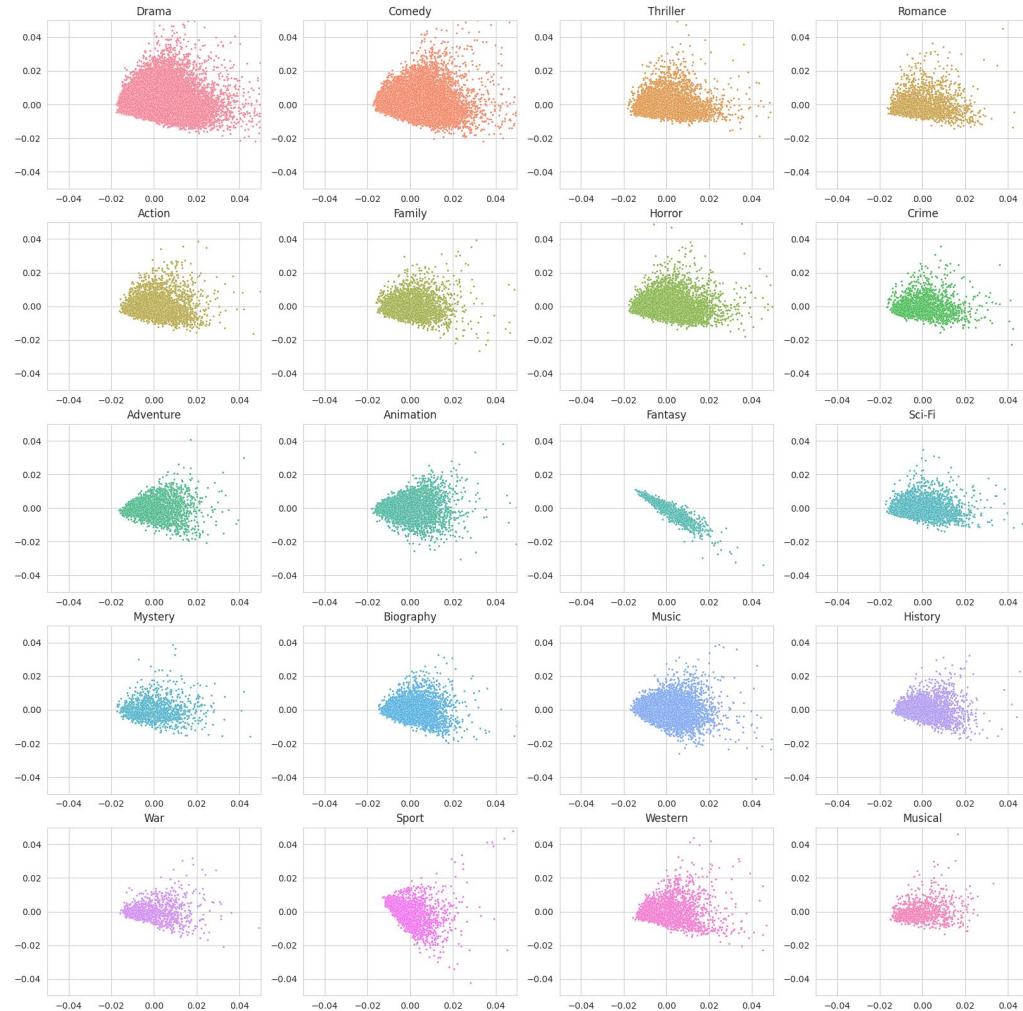
# TF-IDF Weighted Word2Vec Embeddings

Problem with word2vec embeddings: same embeddings for UNK

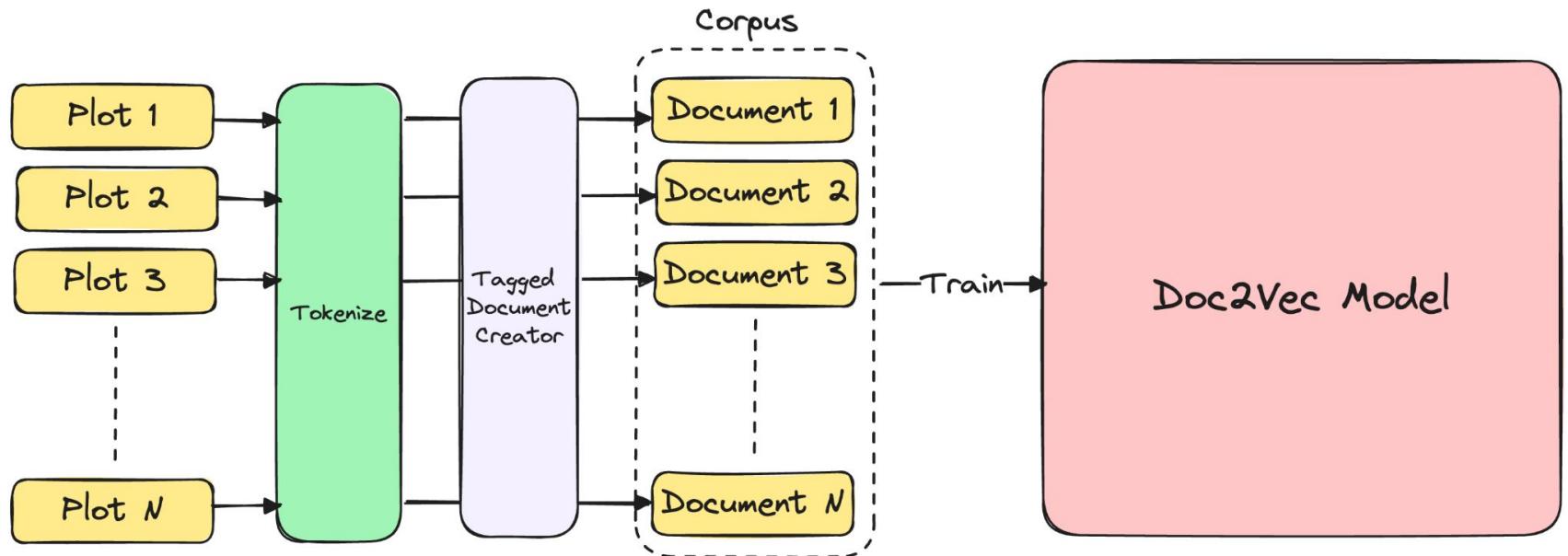
Solved by TF-IDF weighted: a better representation of those UNK tokens



# TF-IDF Weighted Word2Vec Embeddings

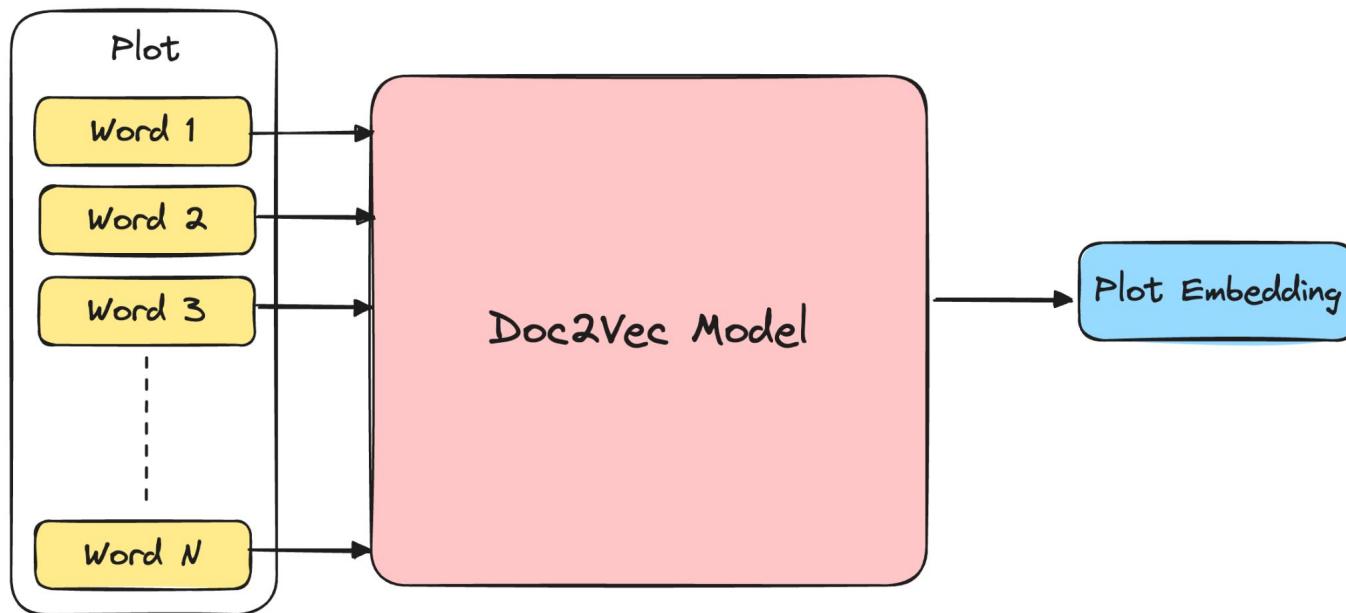


# Doc2Vec Model Training

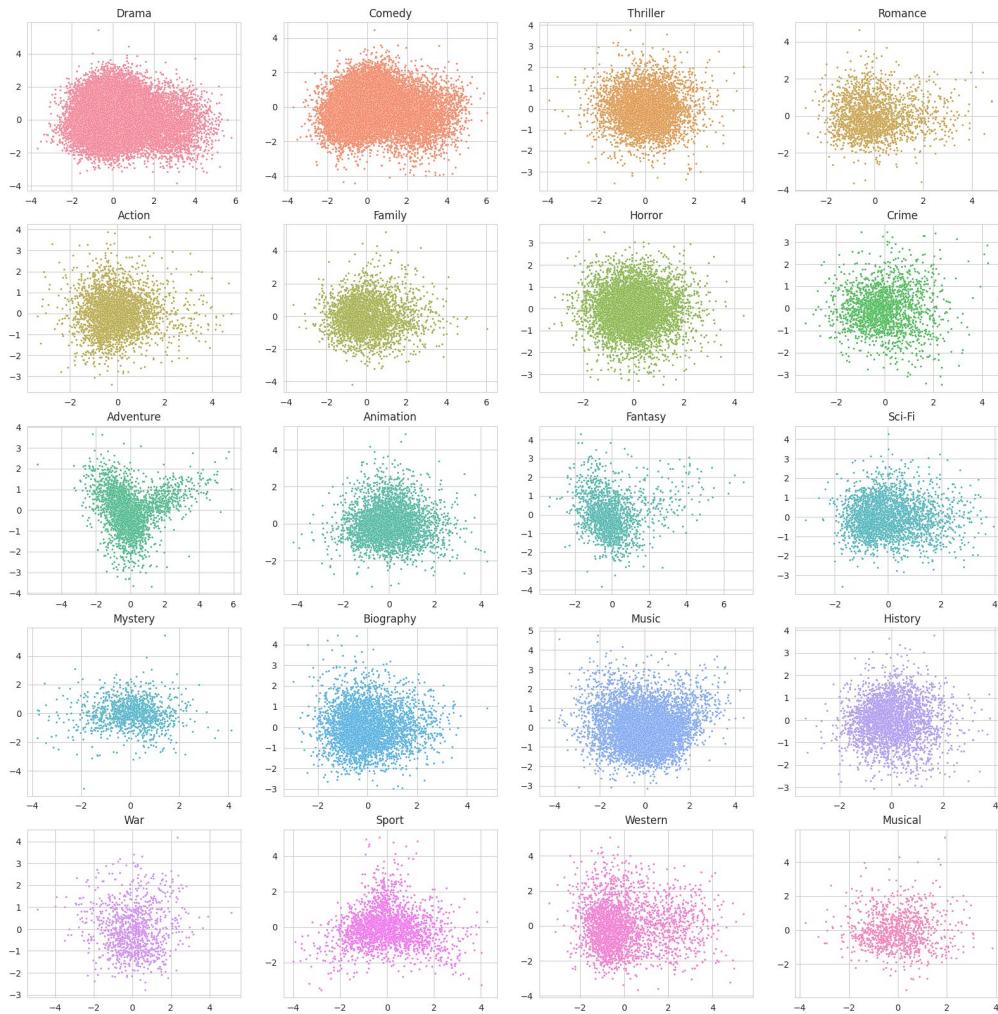


---

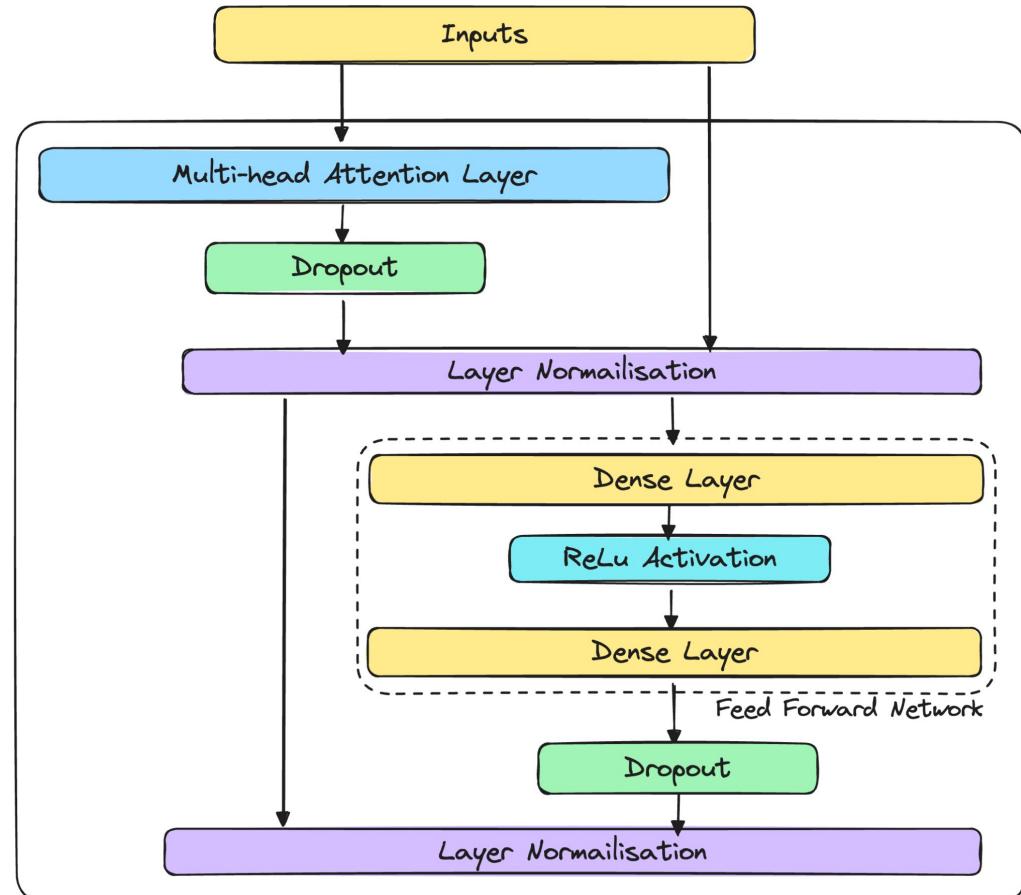
## Doc2Vec Embeddings



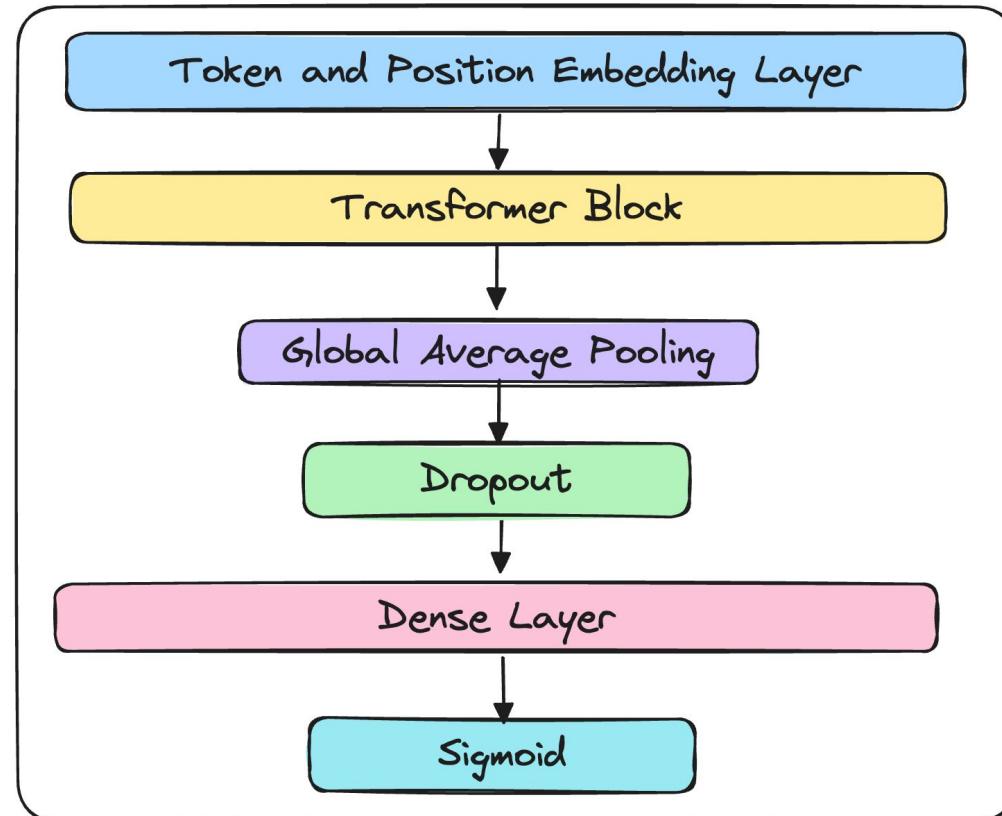
# Doc2Vec Embeddings



# Transformer Block



# Transformer



---

# Results



# Metrics used

**Jaccard Score:** The number of correctly predicted labels divided by the union of predicted and true labels

**Hit Rate:** Proportion of examples where the model predicts at least one correct genre

**F1 Score**

**Precision**

**Recall**

**Accuracy**

**Hamming Score**

---

# Standout Results

# 32.5%

Increase in Jaccard Score for XGBoost with TF-IDF Weighted Word2Vec embeddings over XGBoost with Word2Vec implemented in the paper



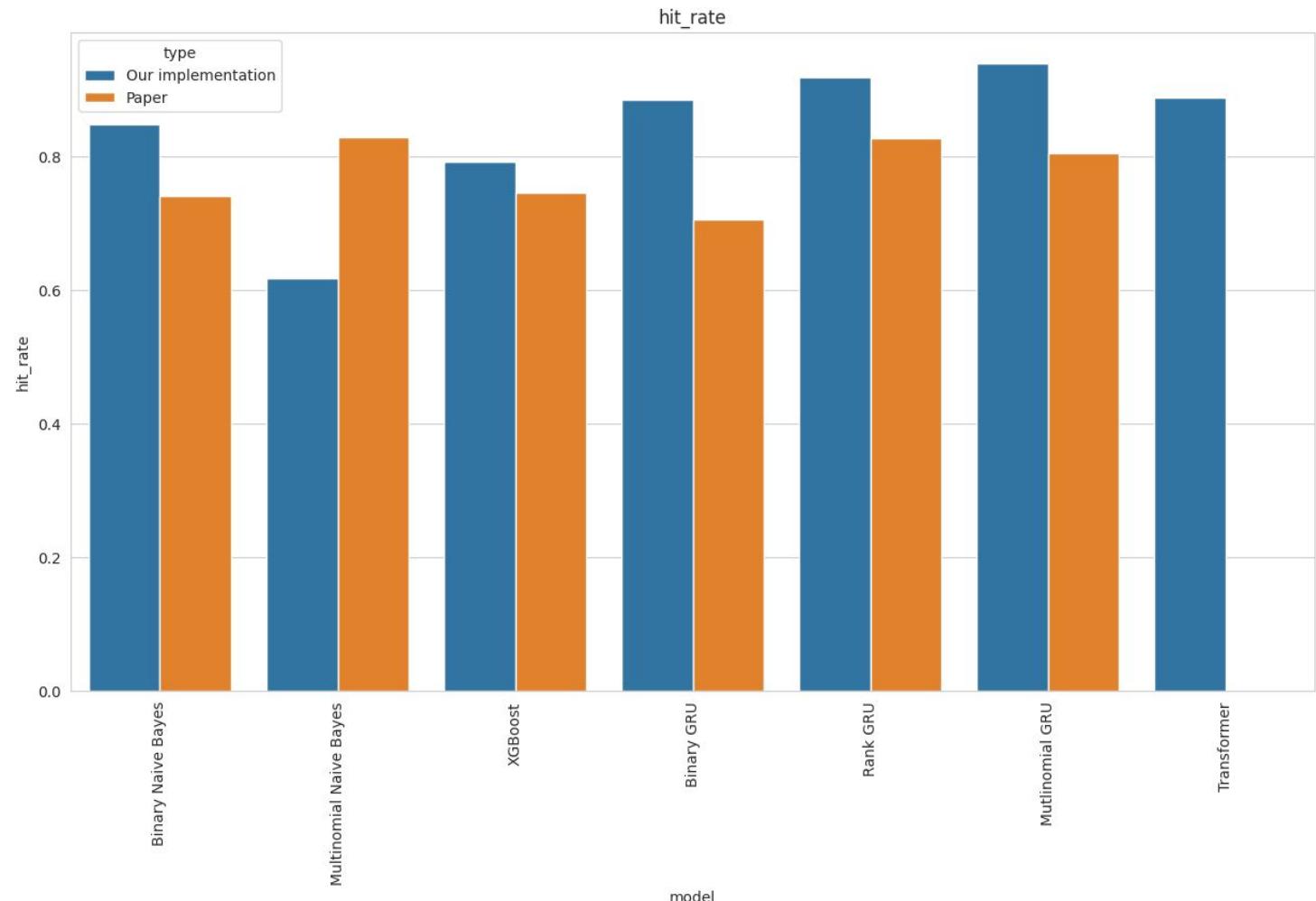
# 50%

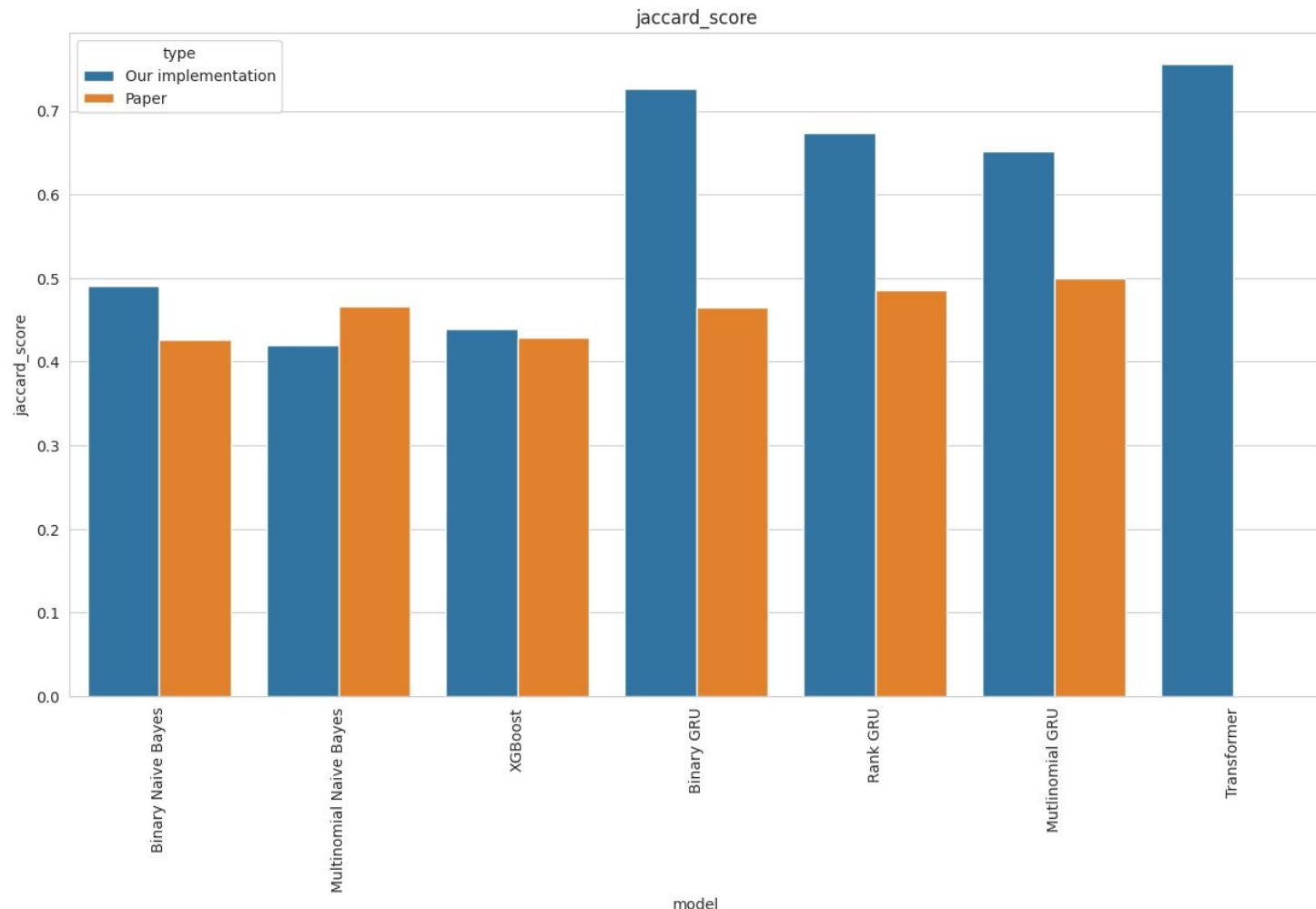
Increase in Jaccard Score for Transformer over Multinomial GRU implemented in the paper

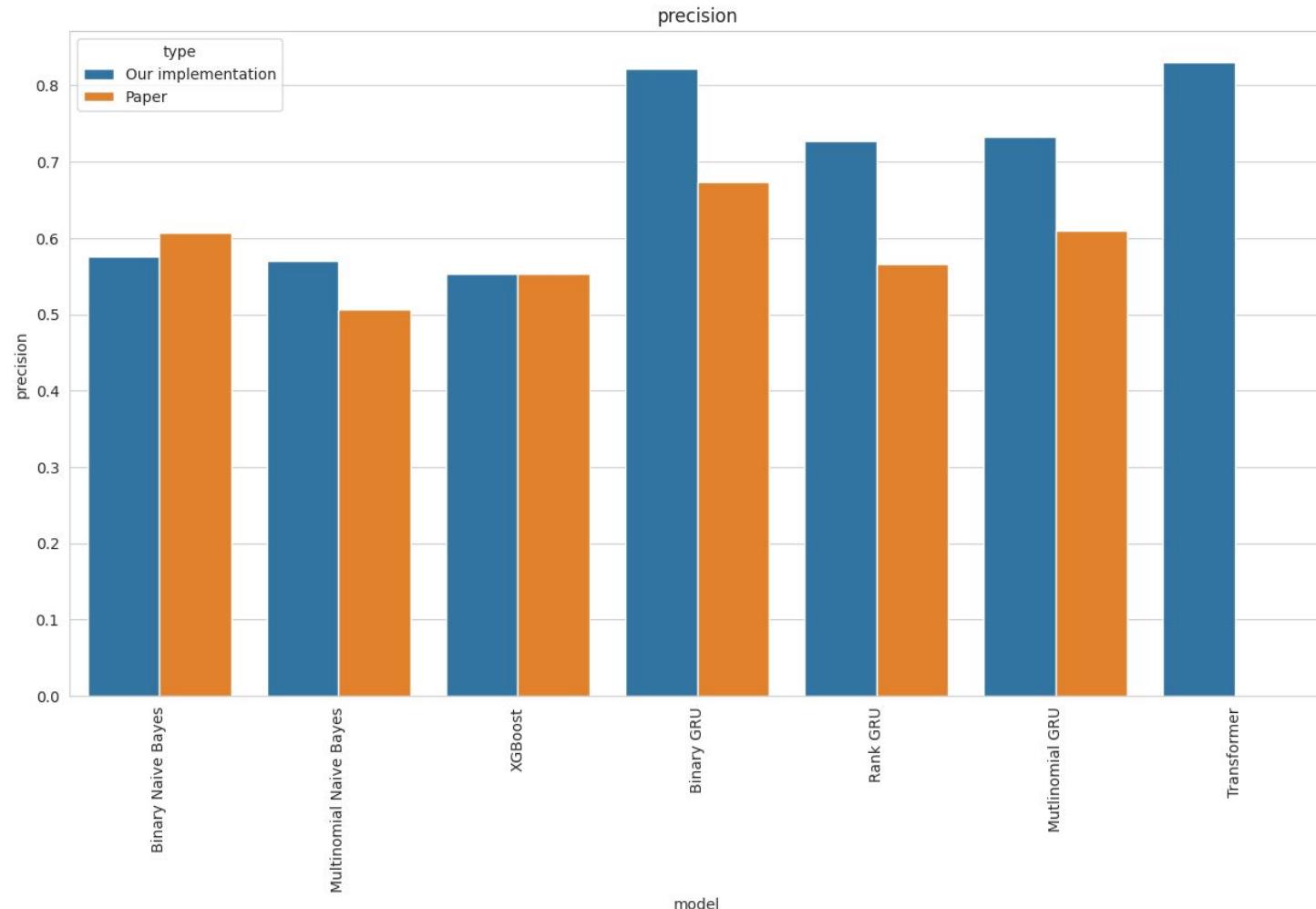


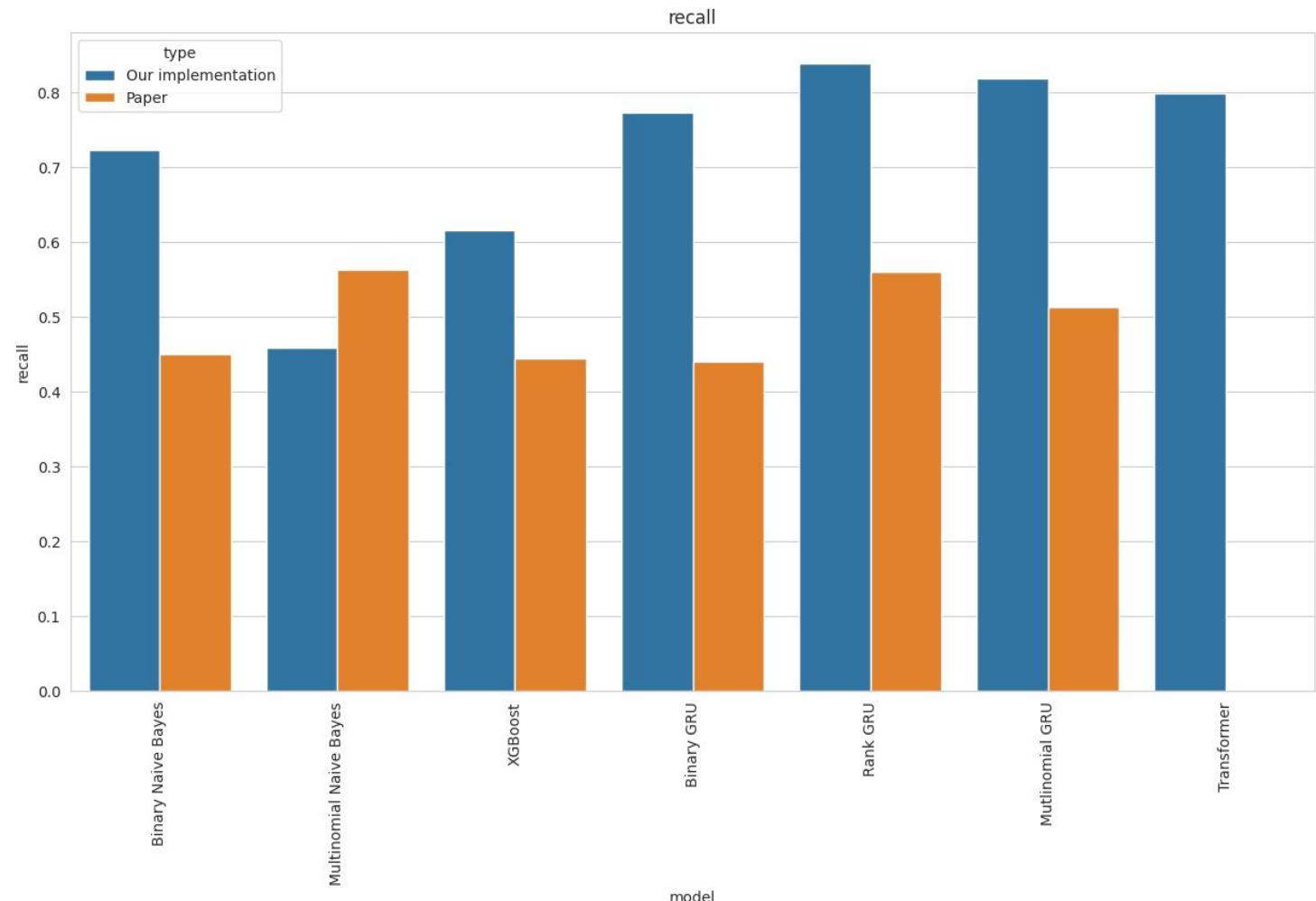
---

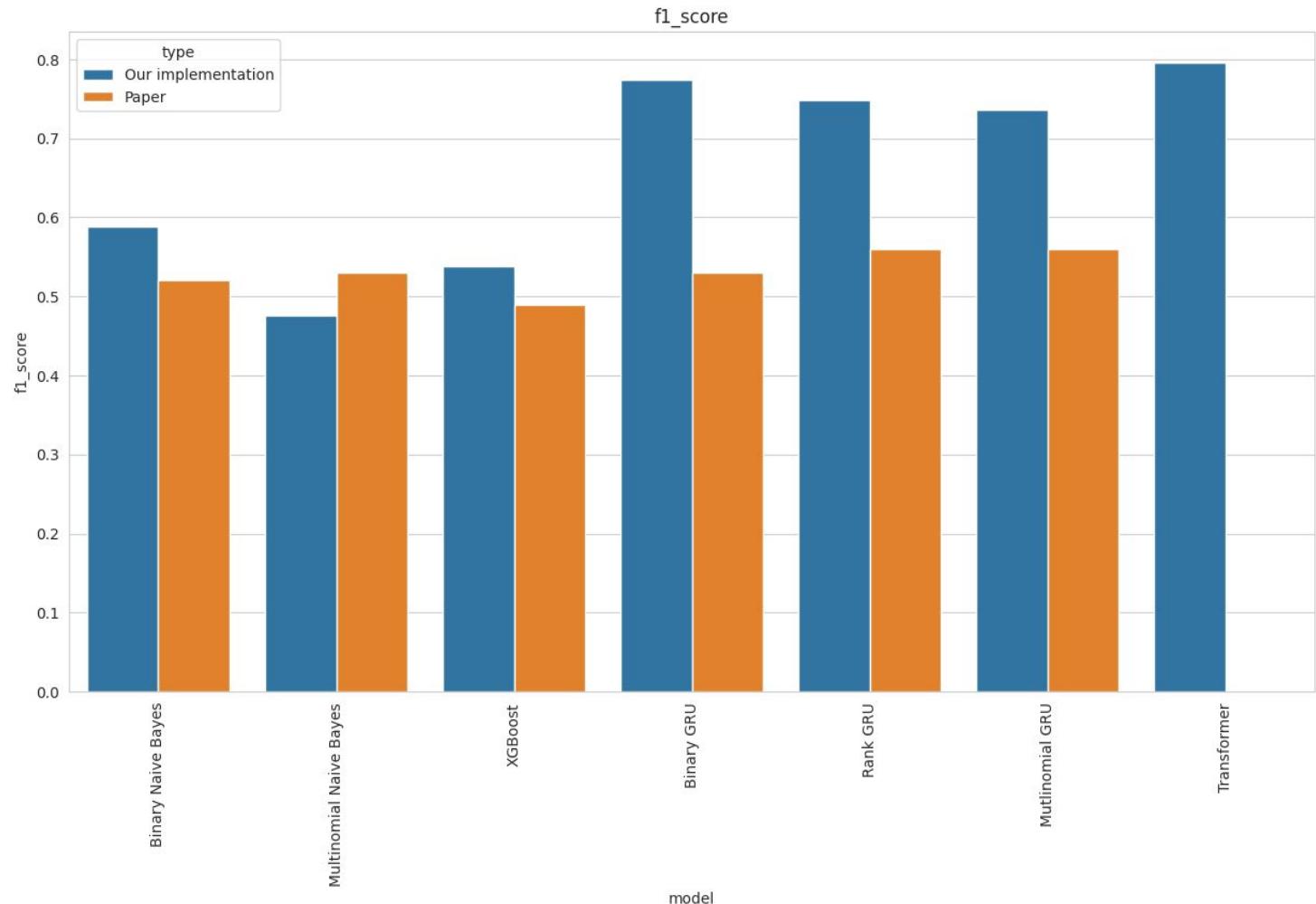
# Standard Results



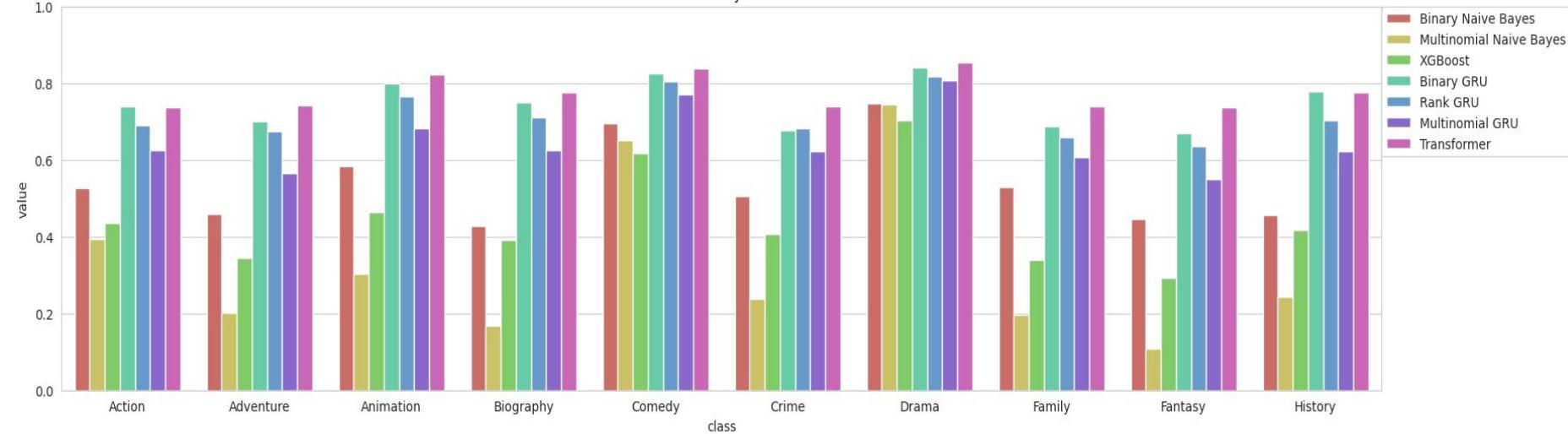




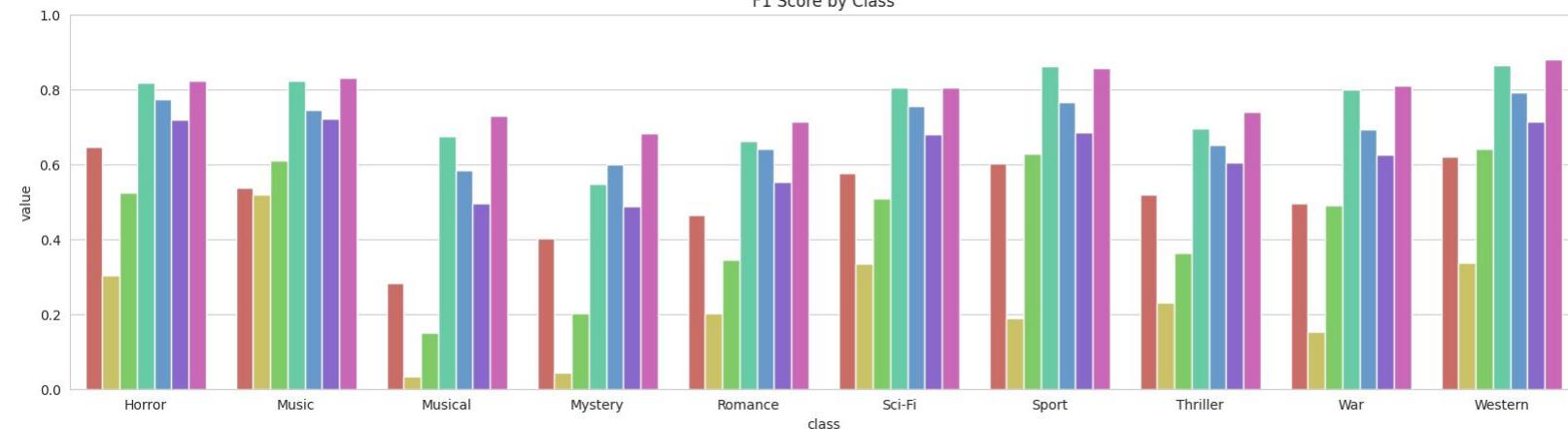




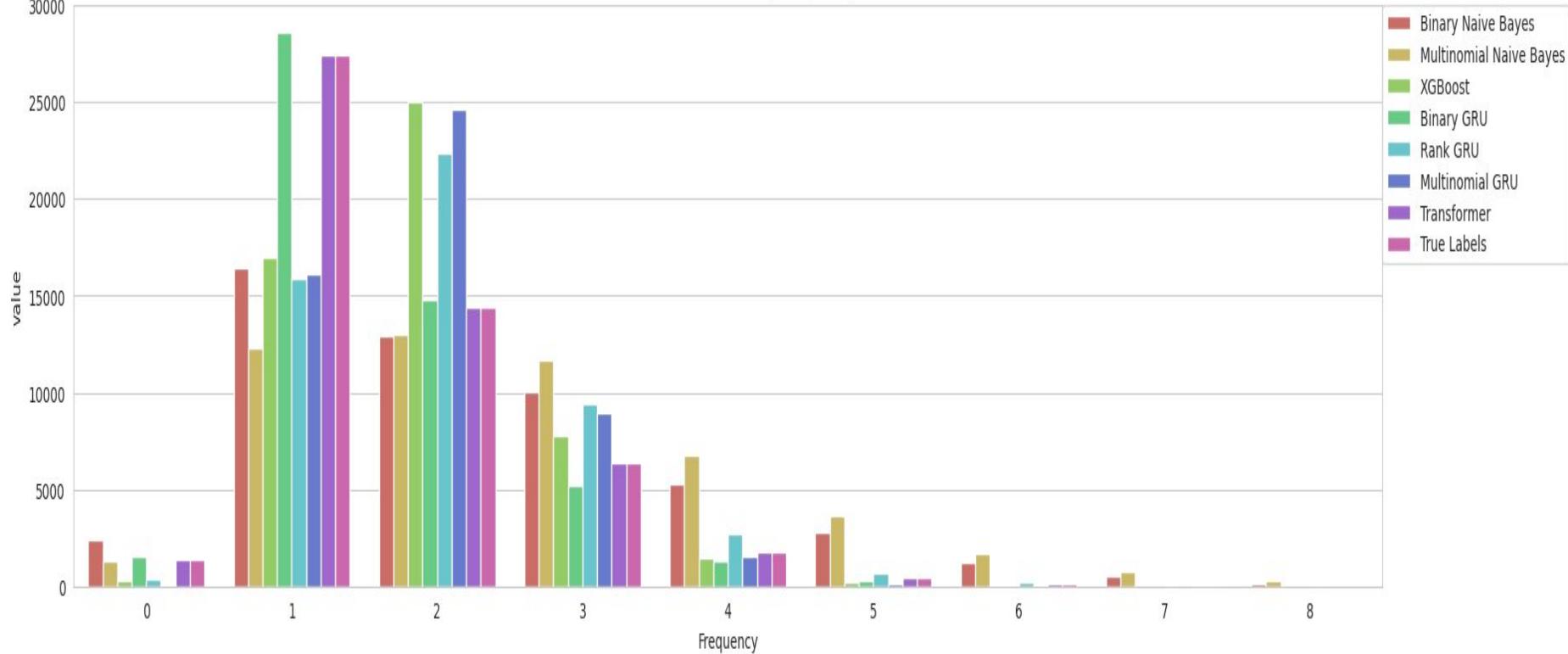
F1 Score by Class



F1 Score by Class



### Number of Predictions by Frequency





**Thank you**