# Benchmark Studio

# LLM MODEL PERFORMANCE COMPARISON ANALYSIS
# Report Milestone 3

**Report by:**

Akshara Sri Lakshmipathy - **akla8196@colorado.edu**

Harish Nandhan Shanmugam - **hash1366@colorado.edu**

Shivaraj Senthil Rajan - **shse1502@colorado.edu**

In this we imported a CSV file (battles_data_cleaned.csv) into a DataFrame and then displayed the first few rows.

| | turn | language | is_code | is_refusal | model | sum_user_tokens | sum_assistant_tokens | context_tokens | dedup_tag_high_freq | dedup_tag_sampled | tstamp_period | information_fulfillment | Math | spec |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | French | False | False | gemma-2-2b-it | 19 | 328 | 19 | False | True | 2024-08-04 | True | False | |
| 1 | 3 | English | False | False | athene-70b-0725 | 51 | 1846 | 1434 | False | True | 2024-08-04 | True | False | |
| 2 | 6 | Italian | False | False | athene-70b-0725 | 472 | 3504 | 3396 | False | True | 2024-08-04 | False | False | |
| 3 | 2 | English | False | False | llama-3-70b-instruct | 84 | 1047 | 580 | False | True | 2024-08-04 | True | False | |
| 4 | 1 | French | False | False | gpt-4-turbo-2024-04-09 | 29 | 556 | 29 | False | True | 2024-08-04 | False | False | |

## HANDLING CLASS IMBALANCE

It calculates the frequency of models in the dataset, filter them and keep only with at least 1000 occurrences and create a new filtered DataFrame containing only these models.

```
model
chatgpt-4o-latest                  9099
gpt-4o-2024-08-06                   5395
llama-3.1-405b-instruct            3557
gemini-1.5-pro-exp-0801            3427
gpt-4o-2024-05-13                   3340
llama-3.1-70b-instruct            3321
claude-3-5-sonnet-20240620        3219
mistral-large-2407                 3113
gemini-1.5-pro-api-0514           3030
gpt-4o-mini-2024-07-18            2535
llama-3.1-8b-instruct             2486
reka-core-20240722                2465
reka-flash-20240722               2334
athene-70b-0725                    2264
gpt-4-turbo-2024-04-09            2047
claude-3-opus-20240229            1939
gemini-1.5-flash-api-0514         1927
gemma-2-27b-it                      1844
deepseek-v2-api-0628              1777
gemma-2-2b-it                      1600
llama-3-70b-instruct              1089
gpt-4-0125-preview                1002
gpt-4-1106-preview                  991
deepseek-coder-v2-0724             987
...
phi-3-mini-4k-instruct-june-2024   628
gemini-advanced-0514               219
mixtral-8x7b-instruct-v0.1         188
Name: count, dtype: int64
```

It calculates the value counts of the 'model' column in the filtered data and prints the number of models.

22

To balance the dataset we can control the number of samples for each unique model. It filters models that have at least 1000 samples, and then resamples each model to ensure an equal representation of samples either by undersampling or oversampling. The code first saves the balanced dataset (balanced_data) to a CSV file and printed the first few rows of the data are displayed to confirm the dataset's successful import in resampled_data.head()

| | turn | language | is_code | is_refusal | model | sum_user_tokens | sum_assistant_tokens | context_tokens | dedup_tag_high_freq | dedup_tag_sampled | tstamp_period | information_fulfillment | Math | spec |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | Persian | False | False | chatgpt-4o-latest | 92 | 227 | 92 | False | True | 2024-08-09 | | False | False |
| 1 | 2 | English | False | True | chatgpt-4o-latest | 1880 | 2504 | 2877 | False | True | 2024-08-08 | | True | False |
| 2 | 1 | English | False | False | chatgpt-4o-latest | 8 | 806 | 8 | False | True | 2024-08-10 | | False | True |
| 3 | 1 | English | False | False | chatgpt-4o-latest | 7 | 78 | 7 | False | True | 2024-08-08 | | False | False |
| 4 | 1 | English | False | False | chatgpt-4o-latest | 36 | 83 | 36 | False | True | 2024-08-11 | | True | False |

After this the code initializes and applies a LabelEncoder to the 'language' column from categorical text to numeric values. A LabelBinarizer is used to binarize the target labels, which helps transform the target variable into a format compatible with the multiclass classification. And then code splits the dataset into training and testing sets (80% training, 20% testing).

## MODEL IMPLEMENTATION

## 1. RANDOM FOREST

A **Random Forest** binary classifier is trained, resulting in multiple classifiers. After training, the classifiers are used to collect the probability scores of each instance in the test set for each class and the class with the highest probability is then selected for each test instance where it calculates and displays metrics—accuracy, precision, recall, and F1-score—for each of the binary classifiers on the test set

```
Metrics for class athene-70b-0725:
  Accuracy: 0.982840909090909
  Precision: 0.888235294117647
  Recall: 0.727710843373494
  F1-score: 0.8
Metrics for class chatgpt-4o-latest:
  Accuracy: 0.9526136363636364
  Precision: 0.20270270270270271
  Recall: 0.0189873417721519
  F1-score: 0.034722222222222224
Metrics for class claude-3-5-sonnet-20240620:
  Accuracy: 0.98125
  Precision: 0.8844827586206897
  Recall: 0.6610824742268041
  F1-score: 0.7566371681415929
Metrics for class claude-3-opus-20240229:
  Accuracy: 0.9872159090909091
  Precision: 0.9104258443465492
  Recall: 0.7908163265306123
  F1-score: 0.8464163822525598
Metrics for class deepseek-v2-api-0628:
  Accuracy: 0.9882386363636364
  Precision: 0.9019886363636364
  Recall: 0.8214747736093143
  F1-score: 0.8598510494245092
...
  Accuracy: 0.9841477272727273
  Precision: 0.9288079470198676
  Recall: 0.7038895859473023
  F1-score: 0.8008565310492506
```
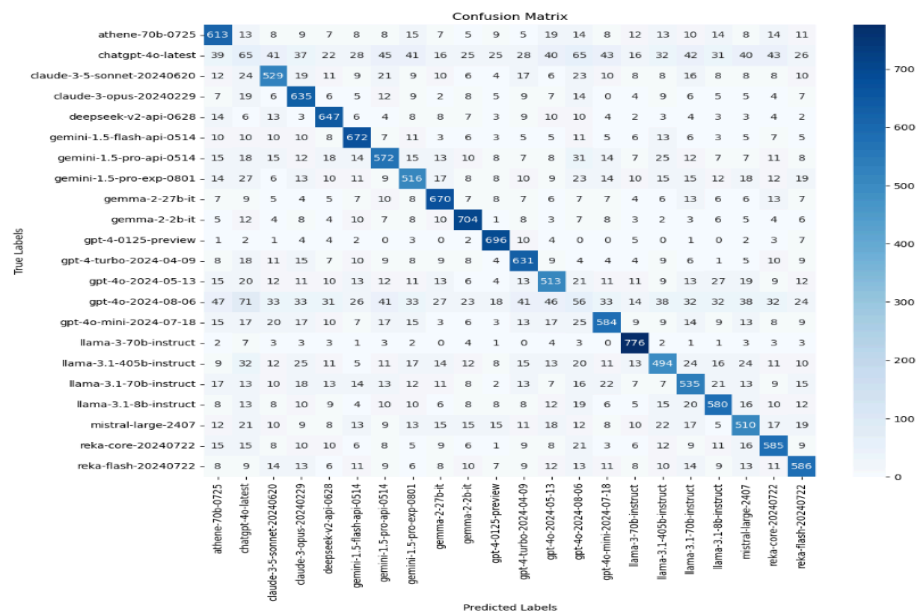
Printed the overall accuracy of the Random Forest classifier across all classes

```
Overall Accuracy of Random Forest Classifier: 0.69
```
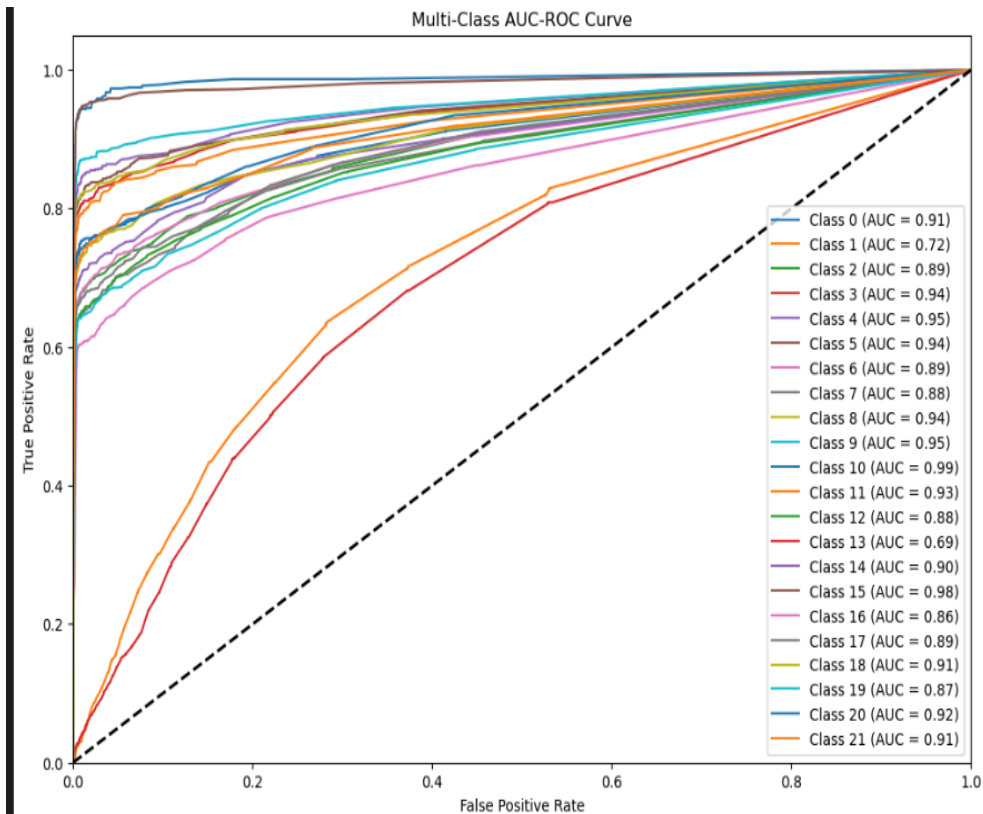
gpt-4-0125-preview and llama-3-70b-instruct showing high precision, recall, and F1-scores, while models like chatgpt-4o-latest underperformed. Most of the models have moderate performance, with scores ranging from 0.60 to 0.80.

```
Classification Report:
                            precision    recall   f1-score    support

          athene-70b-0725        0.69      0.74       0.71        830
        chatgpt-4o-latest        0.15      0.08       0.11        790
claude-3-5-sonnet-20240620        0.68      0.68       0.68        776
   claude-3-opus-20240229        0.69      0.81       0.75        784
     deepseek-v2-api-0628        0.75      0.84       0.79        773
  gemini-1.5-flash-api-0514       0.76      0.82       0.79        821
    gemini-1.5-pro-api-0514       0.68      0.68       0.68        847
    gemini-1.5-pro-exp-0801       0.67      0.65       0.66        796
          gemma-2-27b-it         0.77      0.82       0.79        822
           gemma-2-2b-it         0.79      0.85       0.82        828
         gpt-4-0125-preview       0.83      0.93       0.88        747
    gpt-4-turbo-2024-04-09        0.72      0.79       0.75        799
         gpt-4o-2024-05-13       0.66      0.65       0.66        785
         gpt-4o-2024-08-06       0.13      0.07       0.09        769
     gpt-4o-mini-2024-07-18       0.72      0.70       0.71        840
       llama-3-70b-instruct      0.83      0.94       0.88        825
    llama-3.1-405b-instruct      0.66      0.61       0.63        807
     llama-3.1-70b-instruct      0.66      0.67       0.67        796
      llama-3.1-8b-instruct      0.72      0.73       0.72        797
        mistral-large-2407      0.66      0.65       0.65        789
         reka-core-20240722      0.71      0.75       0.73        782
        reka-flash-20240722      0.72      0.74       0.73        797
...
[ 15  15   8  10  10   6   8   5   9   6   1   9   8  21   3   6  12   9
  11  16 585   9]
[  8   9  14  13   6  11   9   6   8  10   7   9  12  13  11   8  10  14
   9  13  11 586]]
```

The matrix provides a view of how well each model is performing to the others in predicting different classes.



The area under each curve (AUC) is an indicator of how well the model can distinguish between classes, with values closer to 1 indicating better performance.

The Random Forest Classifier undergoes hyperparameter tuning using GridSearchCV to find the optimal parameters, which are used to train a One-vs-Rest classification model. After training, predictions are made by selecting the class with the highest probability for each instance. The evaluation involves calculating overall accuracy, classification report for precision, recall, and F1-score, and visualizing the confusion matrix. Additionally, AUC-ROC curves are plotted.
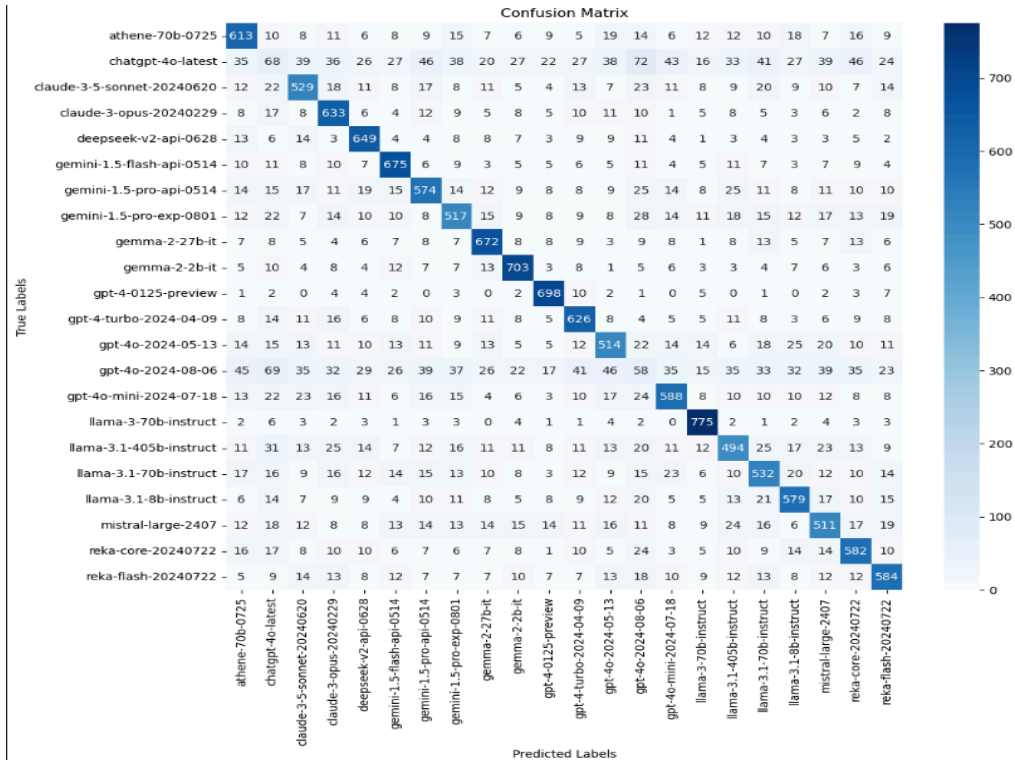
The Random Forest model achieved an overall accuracy of 0.69, with gpt-4-0125-preview and llama-3-70b-instruct performing F1-scores above 0.88. Models like chatgpt-4o-latest had low F1-scores around 0.11.

```
Fitting 3 folds for each of 216 candidates, totalling 648 fits
Best Parameters: {'bootstrap': True, 'max_depth': None, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 300}
Overall Accuracy of Tuned Random Forest Classifier: 0.69
Classification Report:
                           precision    recall  f1-score   support

        athene-70b-0725        0.70      0.74      0.72       830
        chatgpt-4o-latest      0.16      0.09      0.11       790
 claude-3-5-sonnet-20240620    0.67      0.68      0.68       776
   claude-3-opus-20240229      0.70      0.81      0.75       784
     deepseek-v2-api-0628      0.75      0.84      0.79       773
 gemini-1.5-flash-api-0514     0.77      0.82      0.79       821
  gemini-1.5-pro-api-0514      0.69      0.68      0.68       847
  gemini-1.5-pro-exp-0801      0.67      0.65      0.66       796
        gemma-2-27b-it         0.77      0.82      0.79       822
         gemma-2-2b-it         0.79      0.85      0.82       828
      gpt-4-0125-preview       0.83      0.93      0.88       747
    gpt-4-turbo-2024-04-09     0.72      0.78      0.75       799
       gpt-4o-2024-05-13       0.67      0.65      0.66       785
       gpt-4o-2024-08-06       0.14      0.08      0.10       769
    gpt-4o-mini-2024-07-18     0.72      0.70      0.71       840
      llama-3-70b-instruct     0.83      0.94      0.88       825
  llama-3.1-405b-instruct      0.65      0.61      0.63       807
   llama-3.1-70b-instruct      0.65      0.67      0.66       796
    llama-3.1-8b-instruct      0.71      0.73      0.72       797
...

               accuracy                            0.69     17600
              macro avg        0.67      0.69      0.68     17600
           weighted avg        0.67      0.69      0.68     17600
```
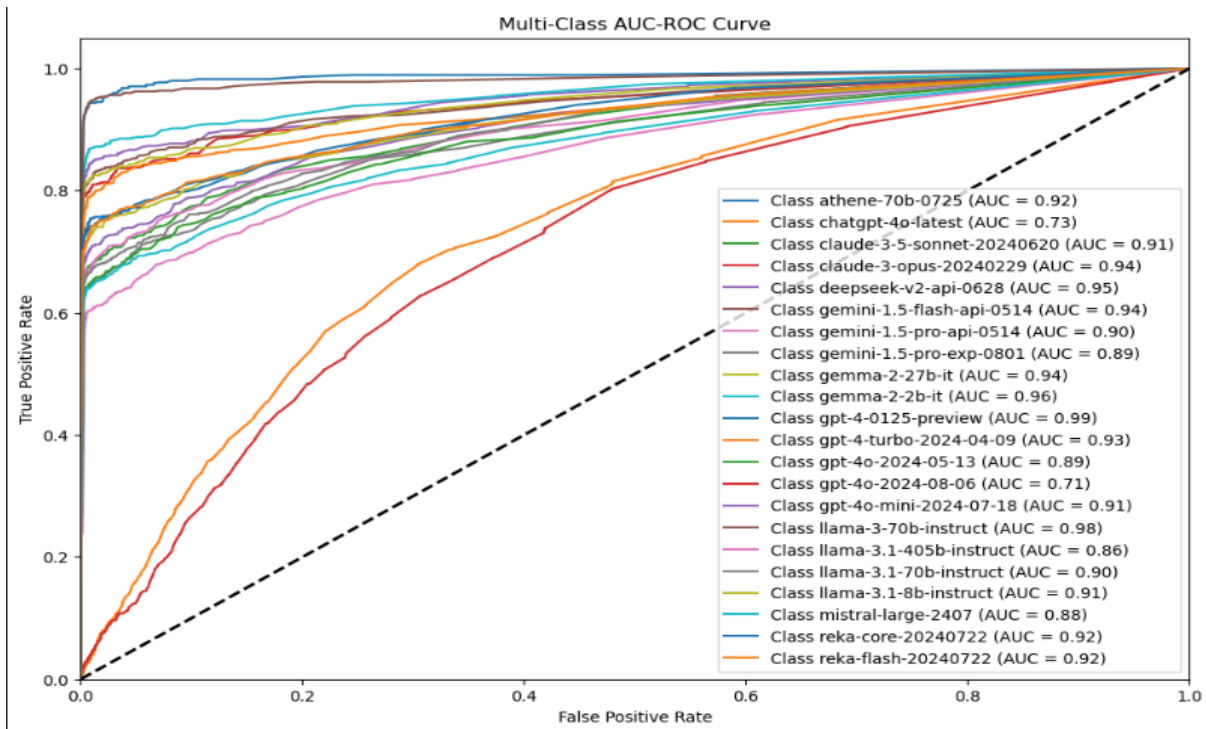
High values on the diagonal elements indicate a high number of correct predictions, while off-diagonal elements show the misclassifications. For example, model 'deepseek-v2-api-0628' has 649 correct predictions for one of its classes, while it has misclassified other classes a few times as shown by the non-diagonal elements.



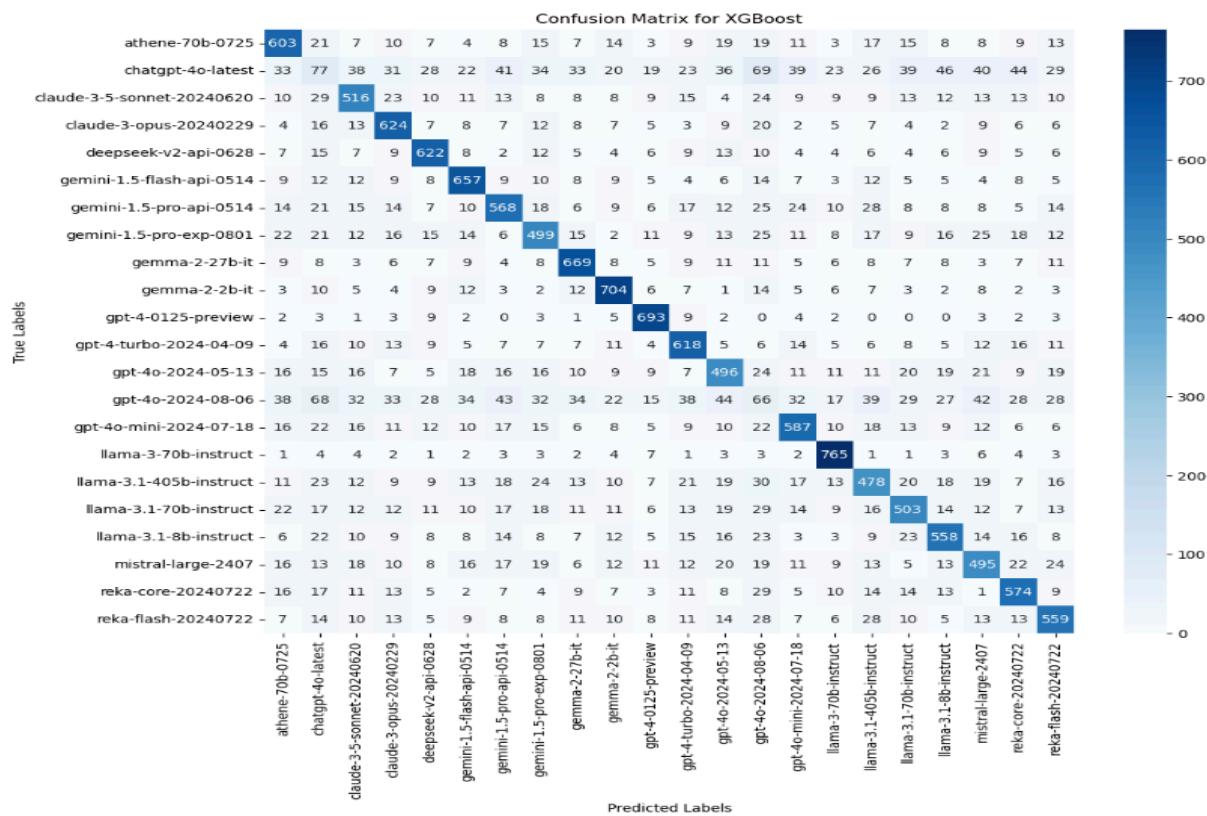Higher AUC values, closer to 1, indicate better model performance.

## 2. XG BOOST

The **XGBoost** classifier achieved an overall accuracy of 0.68.

We have used XGBoost for a multi-class classification using a One-vs-Rest approach and involves the following steps: Data Loading and Preprocessing, Hyperparameter Tuning using RandomizedSearchCV to tune hyperparameters and selecting the best parameters. Then it trains separate One-vs-Rest XGBoost classifiers for each class. Making predictions which involves calculating overall accuracy, classification report for precision, recall, and F1-score, and visualizing the confusion matrix. Additionally, AUC-ROC curves are plotted.
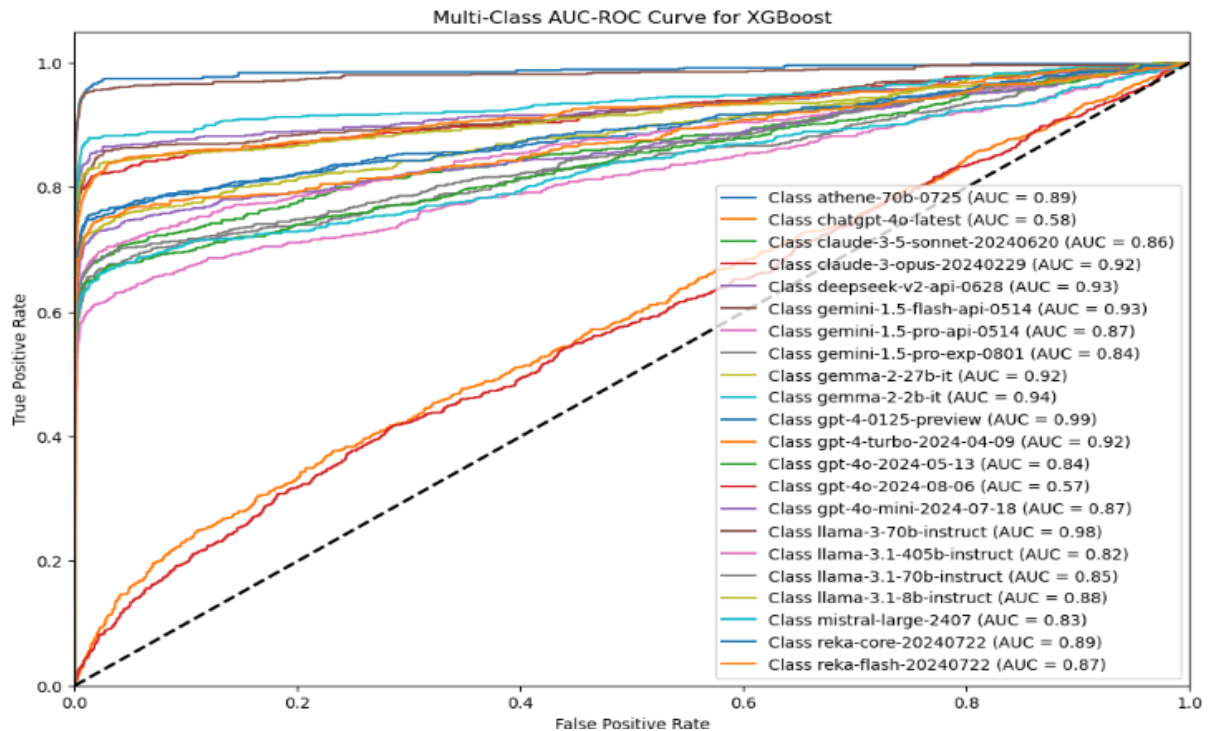
```
Fitting 3 folds for each of 100 candidates, totalling 300 fits
Best Parameters for XGBoost: {'subsample': 1.0, 'scale_pos_weight': 2, 'reg_lambda': 10, 'reg_alpha': 0.1, 'n_estimators': 700, 'min_child_weight': 1, 'max_depth': 15, 'learning_rate':
Overall Accuracy of Tuned XGBoost Classifier: 0.68
Classification Report:
                          precision    recall  f1-score   support

         athene-70b-0725       0.69      0.73      0.71       830
        chatgpt-4o-latest      0.17      0.10      0.12       790
 claude-3-5-sonnet-20240620    0.66      0.66      0.66       776
    claude-3-opus-20240229     0.71      0.80      0.75       784
       deepseek-v2-api-0628    0.75      0.80      0.78       773
  gemini-1.5-flash-api-0514    0.74      0.80      0.77       821
    gemini-1.5-pro-api-0514    0.69      0.67      0.68       847
    gemini-1.5-pro-exp-0801    0.64      0.63      0.64       796
             gemma-2-27b-it    0.75      0.81      0.78       822
              gemma-2-2b-it    0.78      0.85      0.81       828
         gpt-4-0125-preview    0.82      0.93      0.87       747
     gpt-4-turbo-2024-04-09    0.71      0.77      0.74       799
          gpt-4o-2024-05-13    0.64      0.63      0.63       785
          gpt-4o-2024-08-06    0.13      0.09      0.10       769
      gpt-4o-mini-2024-07-18    0.71      0.70      0.71       840
         llama-3-70b-instruct  0.82      0.93      0.87       825
     llama-3.1-405b-instruct   0.62      0.59      0.61       807
      llama-3.1-70b-instruct   0.67      0.63      0.65       796
       llama-3.1-8b-instruct   0.70      0.70      0.70       797
...

                 accuracy                          0.68     17600
                macro avg      0.66      0.68      0.66     17600
             weighted avg      0.66      0.68      0.67     17600
```

The diagonal values, which show the highest numbers (e.g., 603 for the first class), indicate correct predictions.



The curves closer to the top-left corner of the plot indicate better performance, with AUC values closer to 1.0
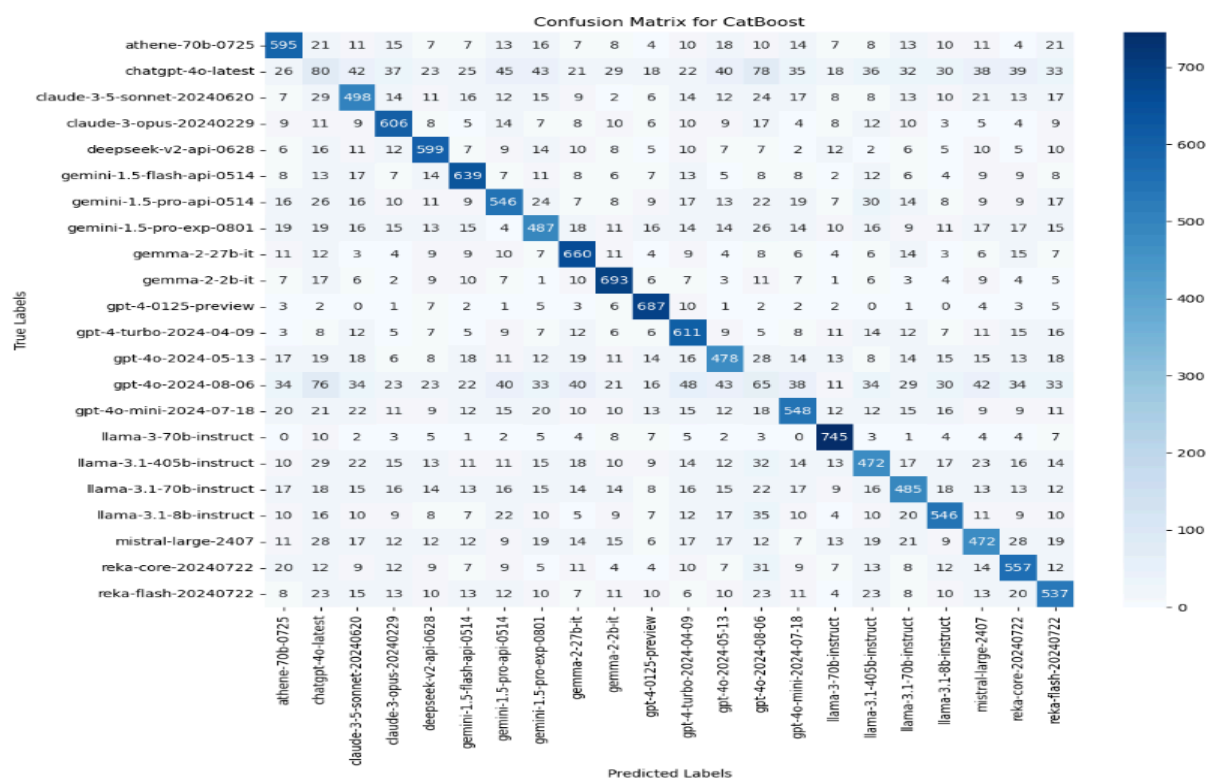
Multi-Class AUC-ROC Curve for XGBoost

## 3. CATBOOST

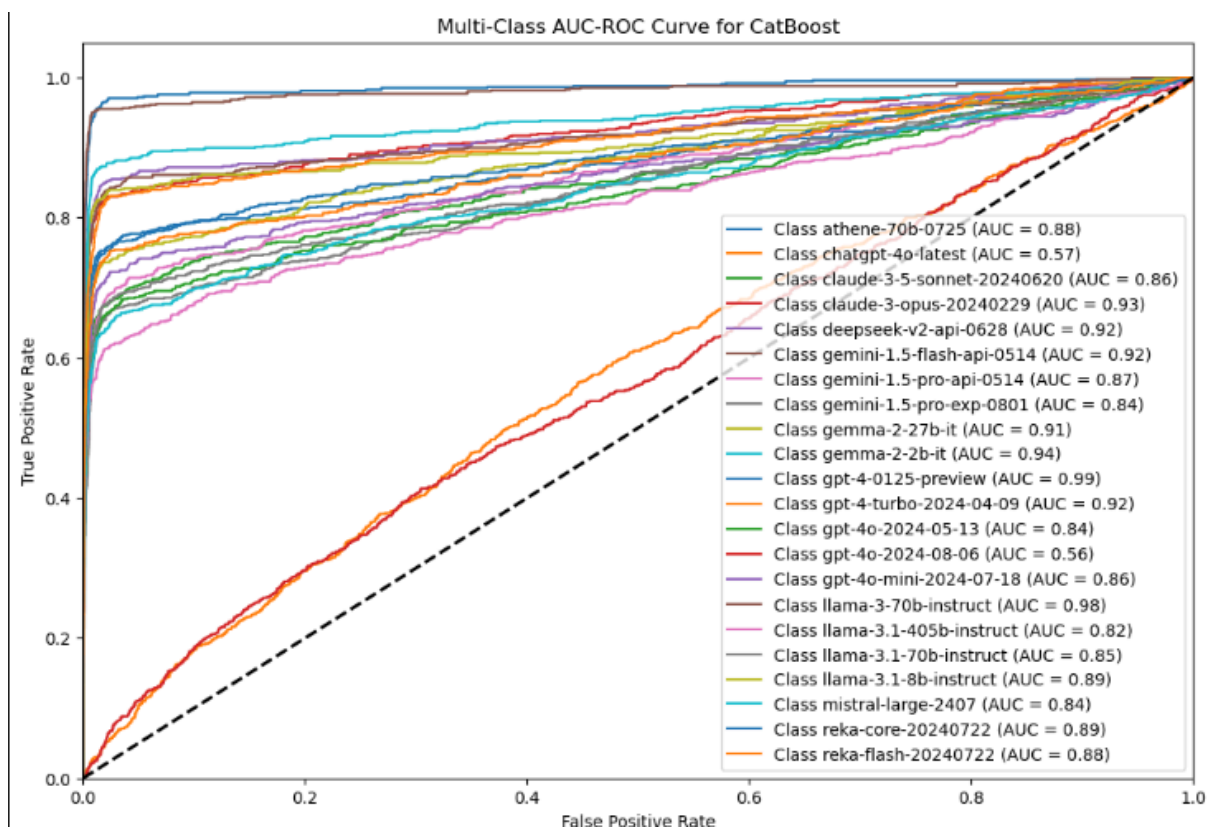The **CatBoost** classifier achieved an overall accuracy of 0.66.

Here we have used CatBoost and involve the following steps: Data Loading and Preprocessing, Hyperparameter Tuning using RandomizedSearchCV to tune hyperparameters and selecting the best parameters. Then it trains separate CatBoost classifiers One-vs-Rest approach. Making predictions which involves calculating overall accuracy, classification report for precision, recall, and F1-score, and visualizing the confusion matrix. Additionally, AUC-ROC curves are plotted.



```
Fitting 3 folds for each of 50 candidates, totalling 150 fits
Best Parameters for CatBoost: {'learning_rate': 0.1388888888888889, 'l2_leaf_reg': 3, 'iterations': 1000, 'depth': 10, 'border_count': 128, 'bagging_temperature': 3}
Overall Accuracy of Tuned CatBoost Classifier: 0.66
Classification Report:
                          precision    recall  f1-score   support

         athene-70b-0725       0.69      0.72      0.71       830
        chatgpt-4o-latest       0.16      0.10      0.12       790
 claude-3-5-sonnet-20240620    0.62      0.64      0.63       776
    claude-3-opus-20240229     0.71      0.77      0.74       784
        deepseek-v2-api-0628   0.72      0.77      0.75       773
    gemini-1.5-flash-api-0514  0.74      0.78      0.76       821
     gemini-1.5-pro-api-0514   0.66      0.64      0.65       847
     gemini-1.5-pro-exp-0801   0.62      0.61      0.62       796
            gemma-2-27b-it     0.72      0.80      0.76       822
             gemma-2-2b-it     0.76      0.84      0.80       828
         gpt-4-0125-preview    0.79      0.92      0.85       747
      gpt-4-turbo-2024-04-09   0.67      0.76      0.72       799
          gpt-4o-2024-05-13    0.64      0.61      0.62       785
          gpt-4o-2024-08-06    0.13      0.08      0.10       769
     gpt-4o-mini-2024-07-18    0.68      0.65      0.67       840
       llama-3-70b-instruct    0.81      0.90      0.85       825
    llama-3.1-405b-instruct    0.62      0.58      0.60       807
     llama-3.1-70b-instruct    0.65      0.61      0.63       796
      llama-3.1-8b-instruct    0.71      0.69      0.70       797
...
                accuracy                           0.66     17600
               macro avg       0.64      0.66      0.65     17600
            weighted avg       0.64      0.66      0.65     17600
```

Confusion Matrix for CatBoost

The Area Under the Curve (AUC) where values close to 1 indicate high ability and values closer to 0.5 suggest no better performance.



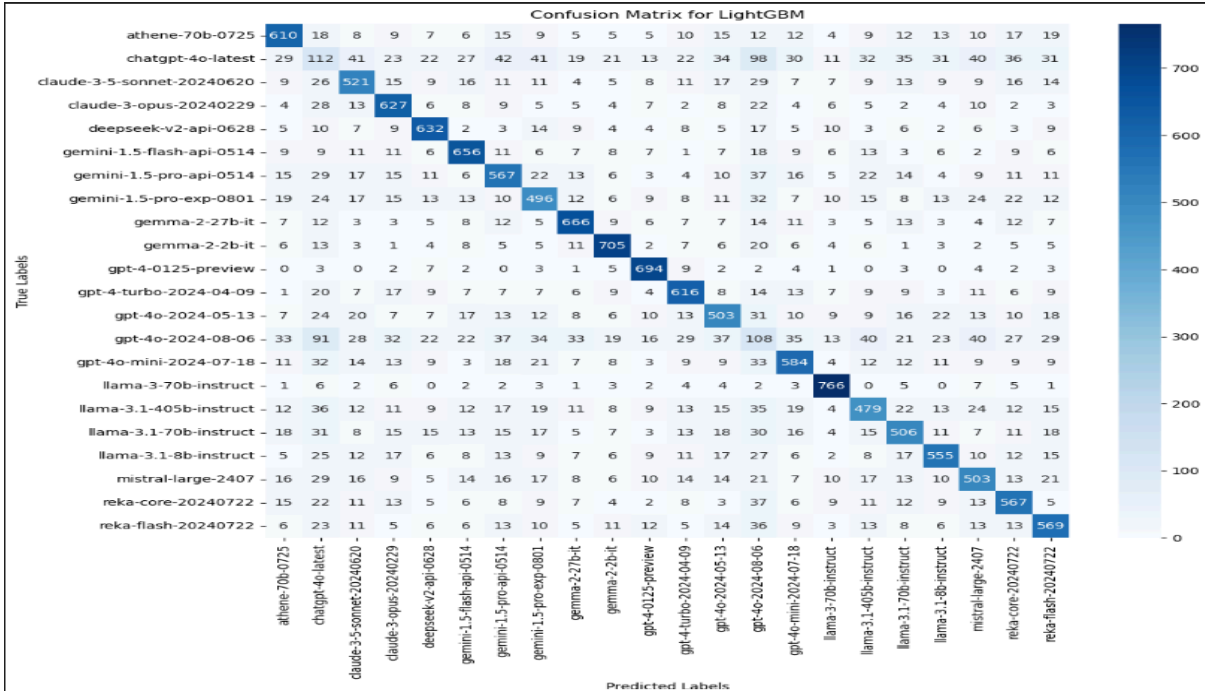Multi-Class AUC-ROC Curve for CatBoost

## 4. LIGHTGBM

The **LightGBM** classifier achieved an overall accuracy of 0.68.

The code uses LightGBM for a multi-class classification through a One-vs-Rest approach and involves the following steps: Data Loading and Preprocessing, Hyperparameter Tuning using RandomizedSearchCV with 150 iterations and 5-fold cross-validation to tune hyperparameters and selecting the best parameters. Then using the One-vs-Rest method, separate LightGBM classifiers and train it. Evaluating the results with confusion matrices, classification reports, and AUC-ROC curves.
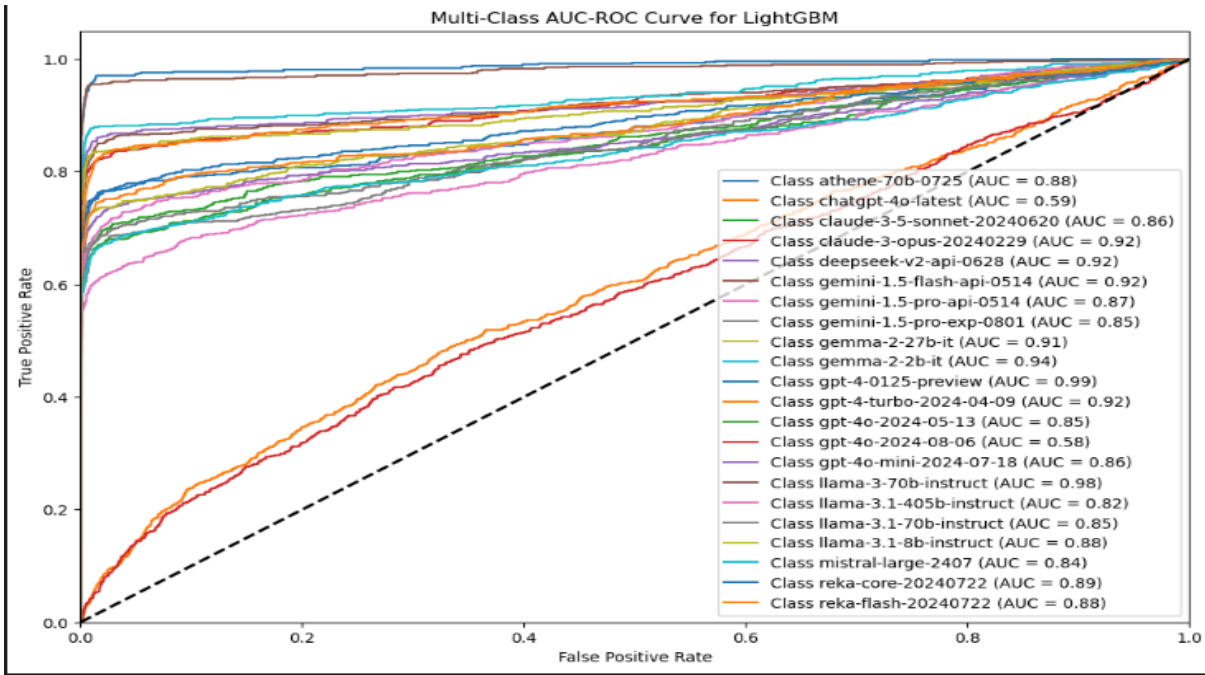
```
Fitting 5 folds for each of 150 candidates, totalling 750 fits
Best Parameters for LightGBM: {'subsample': 0.5526315789473684, 'reg_lambda': 0, 'reg_alpha': 0, 'num_leaves': 150, 'n_estimators': 2000, 'min_split_gain'
Overall Accuracy of Tuned LightGBM Classifier: 0.68
Classification Report:
                            precision    recall  f1-score   support

           athene-70b-0725       0.73      0.73      0.73       830
          chatgpt-4o-latest       0.18      0.14      0.16       790
 claude-3-5-sonnet-20240620       0.67      0.67      0.67       776
      claude-3-opus-20240229       0.72      0.80      0.76       784
         deepseek-v2-api-0628       0.78      0.82      0.80       773
    gemini-1.5-flash-api-0514       0.76      0.80      0.78       821
      gemini-1.5-pro-api-0514       0.67      0.67      0.67       847
      gemini-1.5-pro-exp-0801       0.64      0.62      0.63       796
              gemma-2-27b-it       0.78      0.81      0.80       822
               gemma-2-2b-it       0.82      0.85      0.83       828
          gpt-4-0125-preview       0.83      0.93      0.88       747
       gpt-4-turbo-2024-04-09       0.75      0.77      0.76       799
           gpt-4o-2024-05-13       0.66      0.64      0.65       785
           gpt-4o-2024-08-06       0.16      0.14      0.15       769
       gpt-4o-mini-2024-07-18       0.71      0.70      0.70       840
         llama-3-70b-instruct       0.85      0.93      0.89       825
      llama-3.1-405b-instruct       0.65      0.59      0.62       807
       llama-3.1-70b-instruct       0.67      0.64      0.65       796
        llama-3.1-8b-instruct       0.74      0.70      0.72       797
...
                    accuracy                           0.68     17600
                   macro avg       0.67      0.68      0.68     17600
                weighted avg       0.67      0.68      0.68     17600
```

High diagonal values suggest that the classifier performs well in predicting those classes, whereas high off-diagonal values indicate areas where the classifier confuses one class for another



Values close to 1 indicating excellent ability and values near 0.5 suggesting no better accuracy.
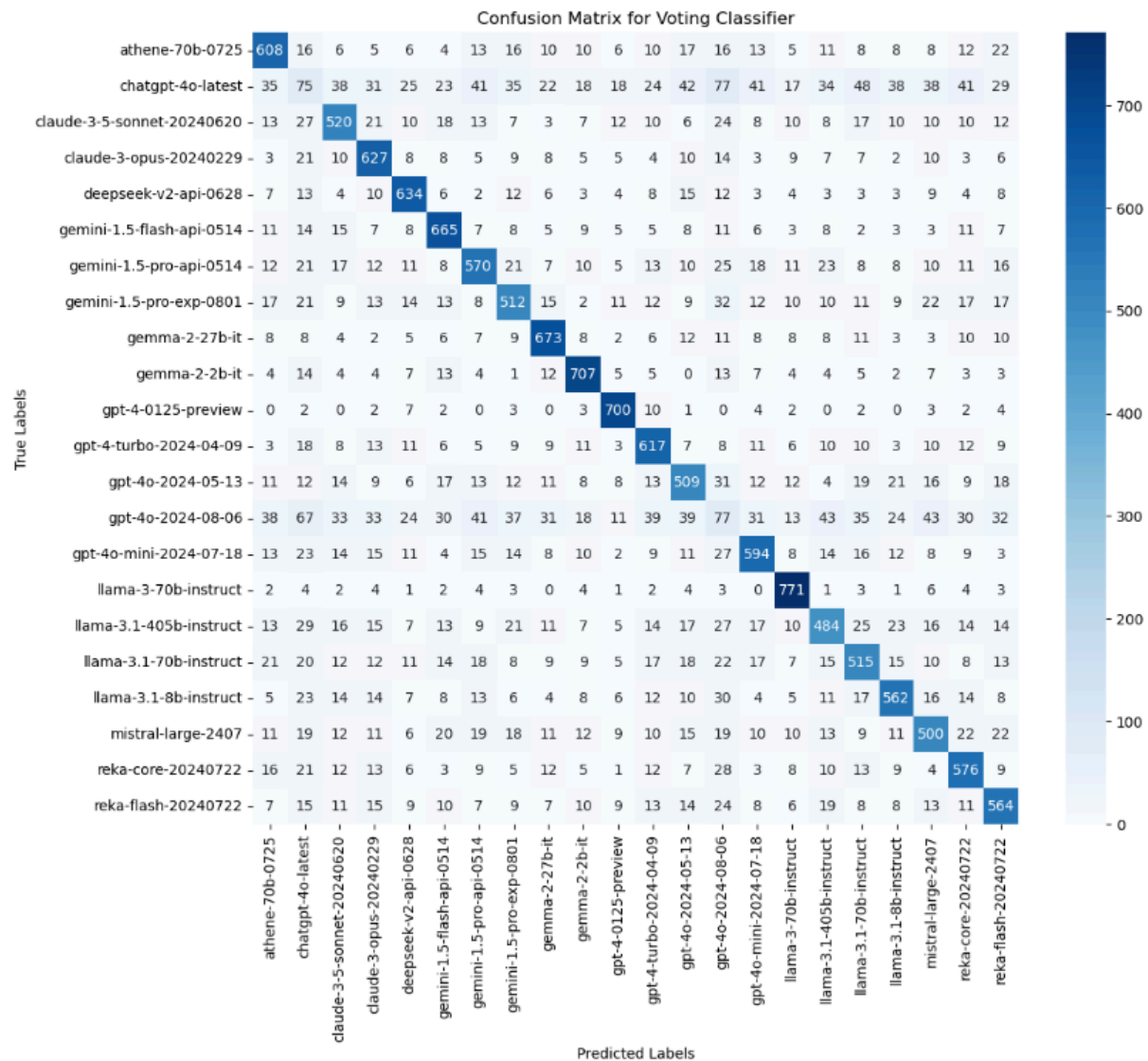
# 5. VOTING

The **Voting** classifier achieved an overall accuracy of 0.69.

We handled a multi-class classification by utilizing a **Voting Classifier** by combining LightGBM, CatBoost, and XGBoost through a One-vs-Rest (OvR) strategy. and involves the following steps: Data Loading and Preprocessing, Model Initialization: where it initializes three models LightGBM, CatBoost, and XGBoost with hyperparameters which constructs a Voting Classifier using soft voting from all three models. Then Train the One-vs-Rest Model with multiple Voting Classifiers in a One-vs-Rest. Evaluating the results with confusion matrices and classification reports.

```
Overall Accuracy of Voting Classifier: 0.69
Classification Report:
                            precision    recall  f1-score   support

          athene-70b-0725       0.71      0.73      0.72       830
         chatgpt-4o-latest       0.16      0.09      0.12       790
  claude-3-5-sonnet-20240620    0.67      0.67      0.67       776
     claude-3-opus-20240229     0.71      0.80      0.75       784
        deepseek-v2-api-0628    0.76      0.82      0.79       773
    gemini-1.5-flash-api-0514   0.74      0.81      0.78       821
      gemini-1.5-pro-api-0514   0.69      0.67      0.68       847
      gemini-1.5-pro-exp-0801   0.66      0.64      0.65       796
             gemma-2-27b-it     0.77      0.82      0.79       822
              gemma-2-2b-it     0.80      0.85      0.83       828
          gpt-4-0125-preview    0.84      0.94      0.89       747
       gpt-4-turbo-2024-04-09   0.71      0.77      0.74       799
          gpt-4o-2024-05-13     0.65      0.65      0.65       785
          gpt-4o-2024-08-06     0.15      0.10      0.12       769
       gpt-4o-mini-2024-07-18   0.72      0.71      0.71       840
        llama-3-70b-instruct    0.82      0.93      0.87       825
      llama-3.1-405b-instruct   0.65      0.60      0.63       807
       llama-3.1-70b-instruct   0.65      0.65      0.65       796
        llama-3.1-8b-instruct   0.73      0.71      0.72       797
          mistral-large-2407    0.65      0.63      0.64       789
          reka-core-20240722    0.69      0.74      0.71       782
...

                 accuracy                           0.69     17600
                macro avg       0.66      0.68      0.67     17600
             weighted avg       0.67      0.69      0.67     17600
```

The numbers on the matrix diagonal (e.g., 608 for the first class) indicate correct predictions for each class, while off-diagonal numbers represent misclassifications.



Confusion Matrix for Voting Classifier

## ACCURACY ACHIEVED BY EACH MODEL

| MODEL NAME | ACCURACY |
|:---:|:---:|
| RANDOM FOREST | 0.69 |
| XG BOOST | 0.68 |
| CATBOOST | 0.66 |
| LIGHTGBM | 0.68 |
| VOTING | 0.69 |