# Phishing Message Detector

As the name implies, the objective is to develop a ML model capable of detecting phishing messages. The model uses only pure text for training (no URLs or any special kind of text).

We use the model, alongside another ML model, that detects phishing URLs, made by our collegue Rafael Pinto, to detect phishing messages and URLs on discord servers. For that, we developed a [discord bot](#).

## Authors

This model was developed and trained by Pompeu Costa and João Mourão.

## Training

To train the model, we used this [dataset from huggingface](#).
Dataset credits go to its author.

## Pre-Processing

Our pre-processing consists of two steps: data pre-processing and feature extraction.

To pre-process, we just take each row of the dataset, transform it into tokens, apply lemmatization and remove stop words.
We then join the pre-processed tokens again to form the pre-processed rows.

For feature extraction, we experimented three different methods: `Bag of Words`, `Fast Text` and `TF-IDF`.
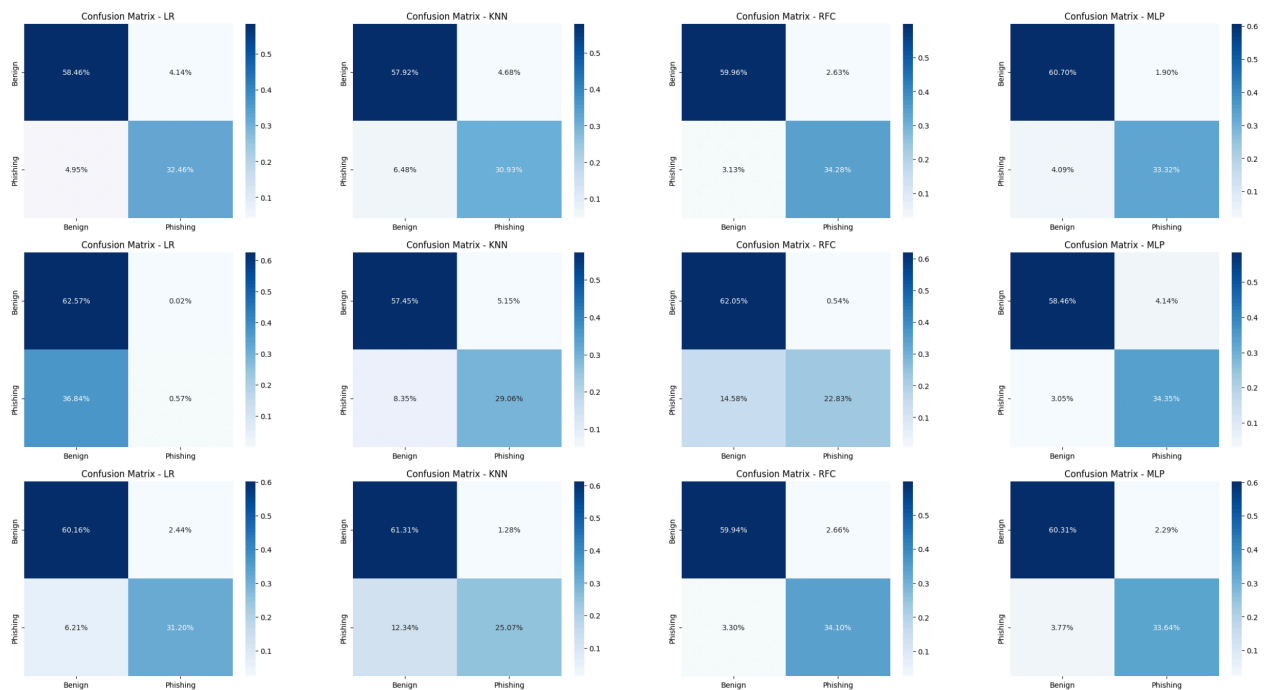In all of the methods, we decided to use a maximum of 512 features.

The results can be seen in their respective md file: [BoW](#), [Fast Text](#) and [TF-IDF](#).
The parameters used for each classifier are described in the next section.

To facilitate the viewing of the results, the image below contains the results of each feature extraction method.
The columns represent the classifier and the rows represent the method (BoW -> Fast Text -> TF-IDF).

Confusion Matrix - LR | Confusion Matrix - KNN | Confusion Matrix - RFC | Confusion Matrix - MLP

Row 1:
- LR: Benign 58.46%, 4.14%; Phishing 4.95%, 32.46%
- KNN: Benign 57.92%, 4.68%; Phishing 6.48%, 30.93%
- RFC: Benign 59.96%, 2.63%; Phishing 3.13%, 34.28%
- MLP: Benign 60.70%, 1.90%; Phishing 4.09%, 33.32%

Row 2:
- LR: Benign 62.57%, 0.02%; Phishing 36.84%, 0.57%
- KNN: Benign 57.45%, 5.15%; Phishing 8.35%, 29.06%
- RFC: Benign 62.05%, 0.54%; Phishing 14.58%, 22.83%
- MLP: Benign 58.46%, 4.14%; Phishing 3.05%, 34.35%

Row 3:
- LR: Benign 60.16%, 2.44%; Phishing 6.21%, 31.20%
- KNN: Benign 61.31%, 1.28%; Phishing 12.34%, 25.07%
- RFC: Benign 59.94%, 2.66%; Phishing 3.30%, 34.10%
- MLP: Benign 60.31%, 2.29%; Phishing 3.77%, 33.64%

# Classifiers

After pre-processing, the model needs to be trained. For that, we experimented with four classifiers: `Logistic Regression`, `K Nearest Neighbours`, `Random Forest` and `Multi-Layer Perceptron`.

To assess the best parameters for each classifier, we performed a `Grid Search` for the following parameters:

```python
params = {
    'LR': {
        'C': [0.01, 0.1, 1, 10],
        'solver': ['lbfgs', 'liblinear','saga','sag'],
        'max_iter': [100, 200, 300]
    },
    'KNN': {
        'n_neighbors': [3, 5, 7, 10],
        'weights': ['uniform', 'distance'],
        'metric': ['euclidean', 'manhattan']
    },
    'RFC': {
        'n_estimators': [50, 100, 150, 200],
        'max_depth': [None, 10, 20, 30],
        'min_samples_split': [2, 5, 10],
        'min_samples_leaf': [1, 2, 4]
    },
    'MLP': {
        'hidden_layer_sizes': [(1,),(100, 100),(128,64,32), (256,128,64)],
        'activation': ['relu', 'tanh', 'identity', 'logistic'],
        'solver': ['adam', 'sgd', 'lbfgs'],
        'learning_rate_init': [0.001, 0.01],
```

```
        }
    }
```

The best parameters for each classifier were the following:

- LR: {'C': 0.1, 'max_iter': 100, 'solver': 'lbfgs'}
- KNN: {'metric': 'euclidean', 'n_neighbors': 5, 'weights': 'distance'}
- RFC: {'max_depth': None, 'min_samples_leaf': 1, 'min_samples_split': 10, 'n_estimators': 200}
- MLP: {'activation': 'relu', 'hidden_layer_sizes': (256, 128, 64), 'learning_rate_init': 0.001,'solver': 'adam'}
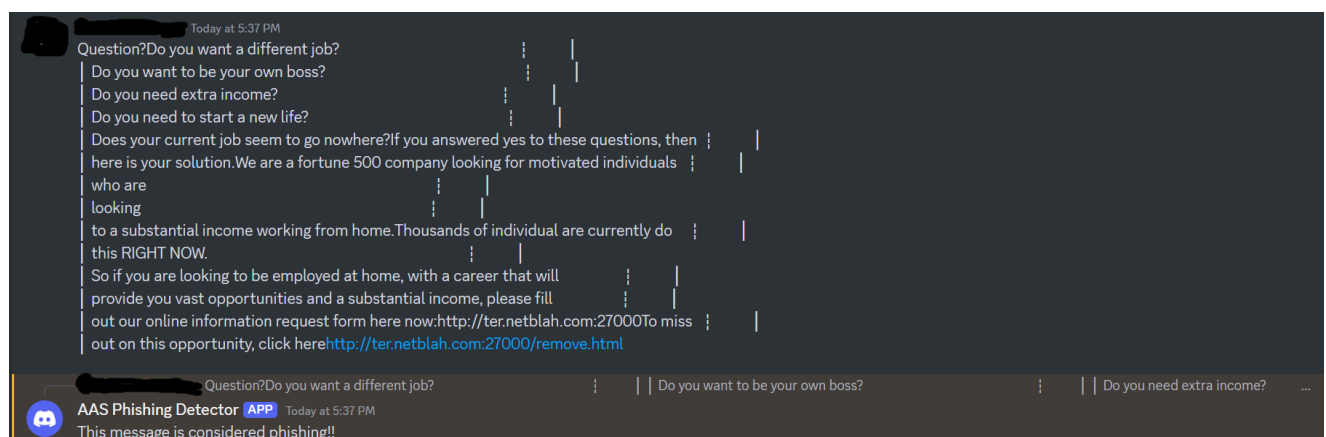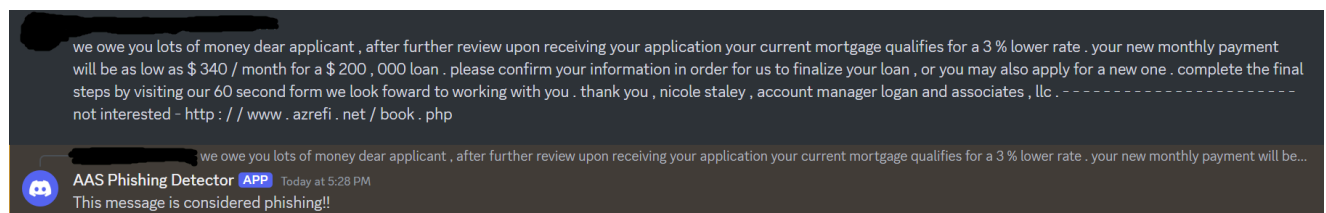
# Results

Looking at the compilation image, its possible to see that the best one is the `Multi-Layer Perceptron` classifier with the `Fast Text` method.
It keeps a good accuracy while having a low rate of false positives.

Having a benign message that is considered phishing will not cause any harm, however, having a phishing message that is considered benign can cause harm to the users.
Thus, we decided to sacrifice a bit of accuracy in favour of lower false positive rates.
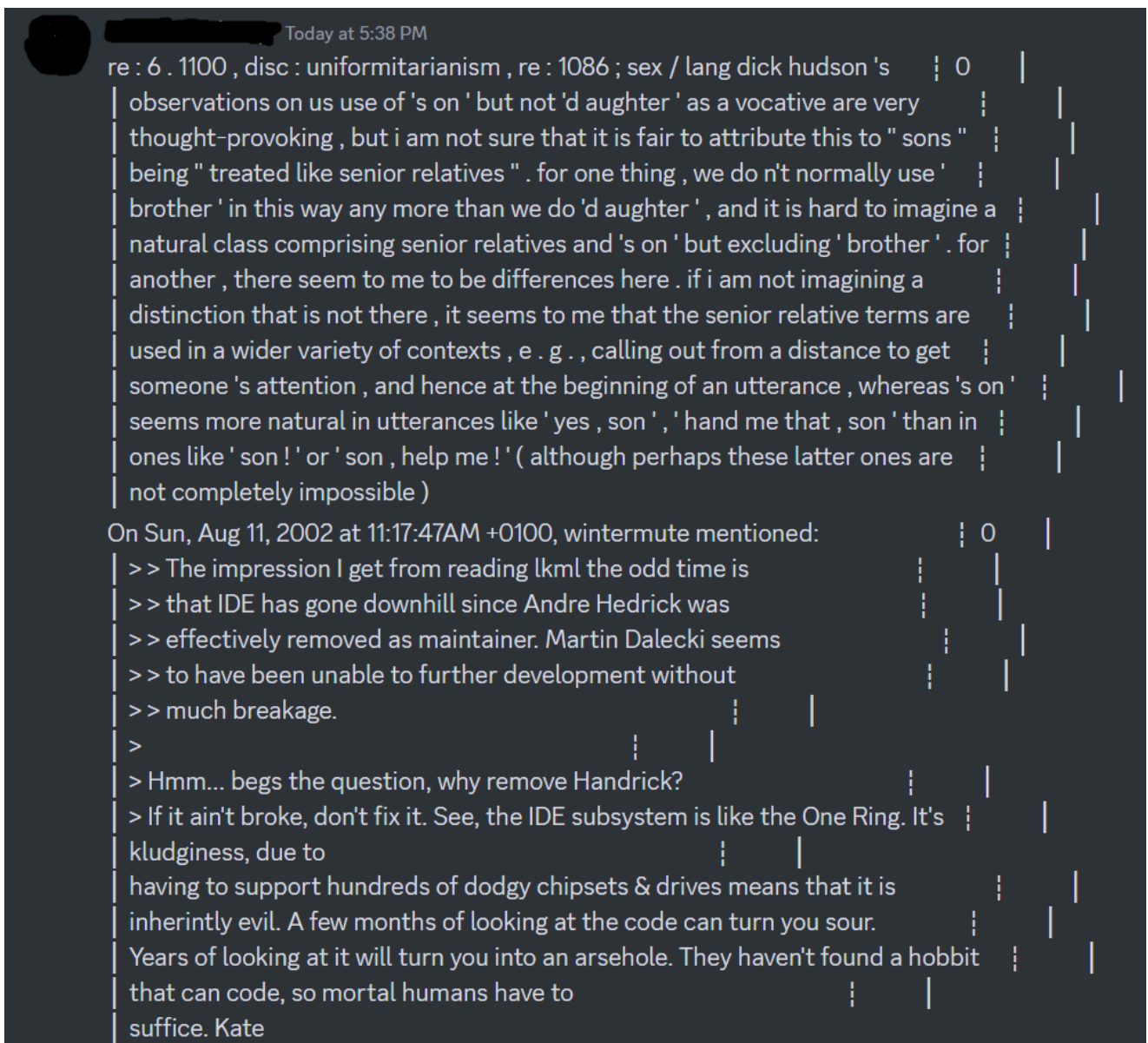
# Running the model in the discord bot

The bot, using this model, was capable of identifying phishing messages.
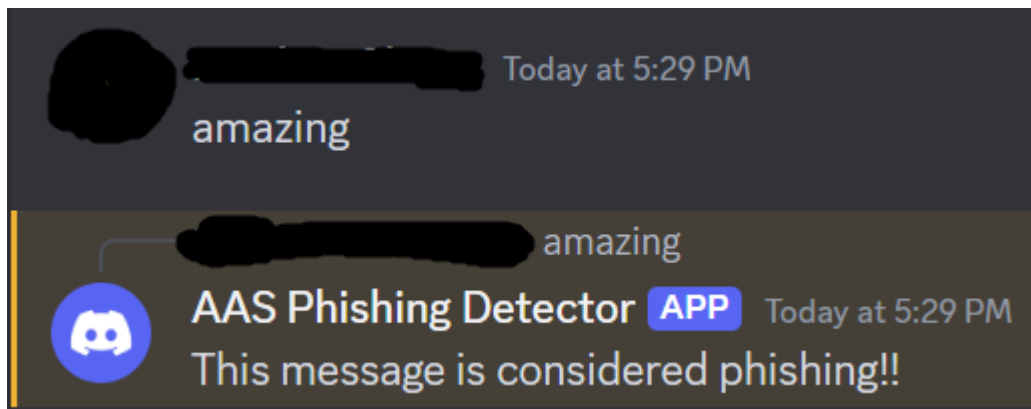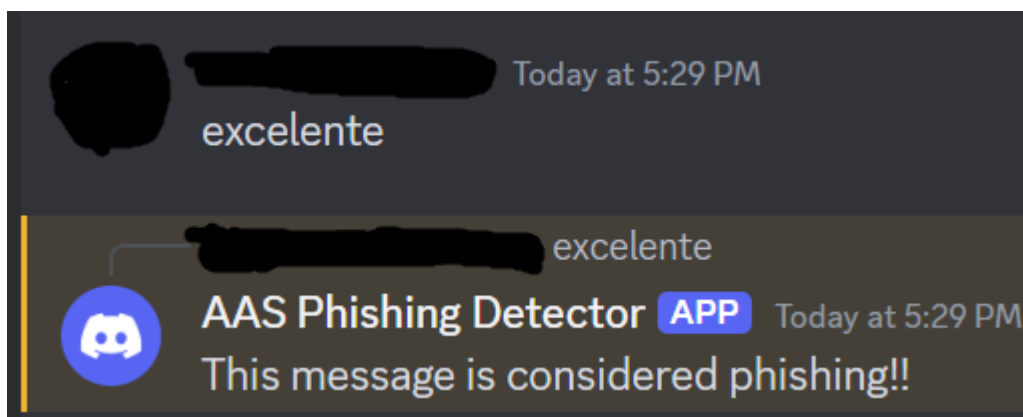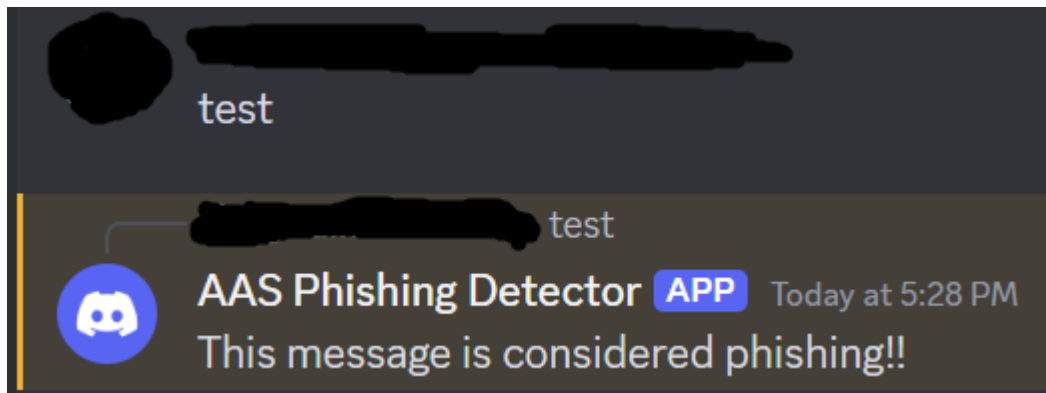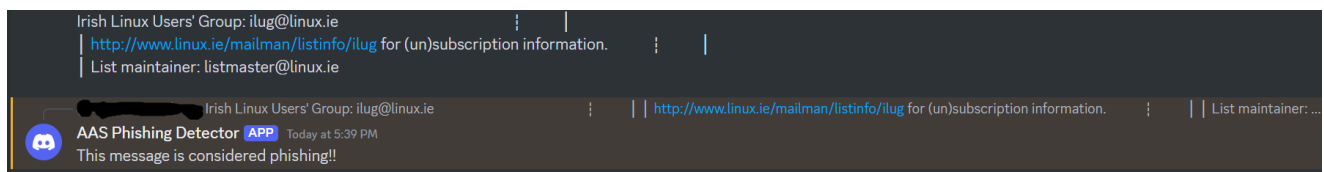
It also successfully identified benign messages (the bot doesn't reply to the message when its considered benign).

However, its not perfect. It considered some benign messages as phishing messages.

We suspect that its due to the `Fast Text` model not being trained well enough.

Nevertheless, the model and the bot worked as intended.
Albeit not perfect, it detects phishing messages.