

# SSH Shell Attacks - Appendix

ANDREA BOTTICELLA\*, Politecnico di Torino, Italy  
ELIA INNOCENTI\*, Politecnico di Torino, Italy  
RENATO MIGNONE\*, Politecnico di Torino, Italy  
SIMONE ROMANO\*, Politecnico di Torino, Italy

## CONTENTS

Contents	1
A DATA EXPLORATION AND PRE-PROCESSING	1
B SUPERVISED LEARNING - CLASSIFICATION	4
C UNSUPERVISED LEARNING - CLUSTERING	7
D LANGUAGE MODEL EXPLORATION	9

### A DATA EXPLORATION AND PRE-PROCESSING

This appendix contains additional plots and visualizations related to the data exploration and pre-processing phase of the analysis. The figures are grouped into subsections based on their thematic relevance.

#### A.1 Temporal Analysis of Attacks

##### A.1.1 Attack Frequency by Hour

The plot shows the distribution of SSH attacks across different hours of the day. It reveals specific hours when attack activity peaks, which could indicate targeted times for malicious activities. Understanding these patterns can help in implementing time-based security measures.

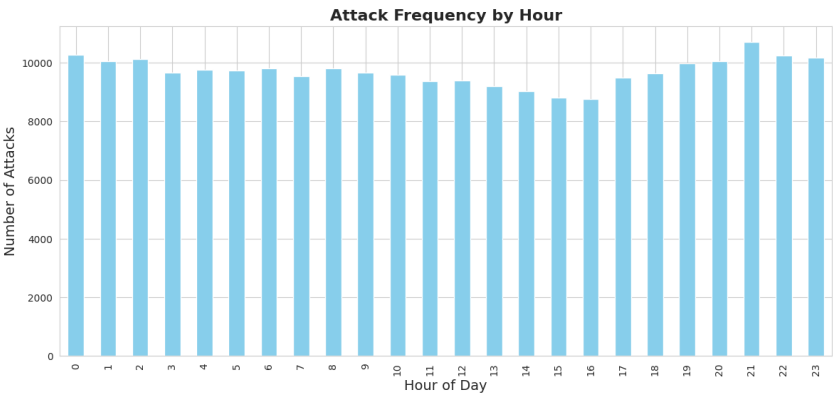


Fig. 1. Distribution of SSH attacks by hour of the day - The plot highlights peak hours during which attacks are most frequent

\*The authors collaborated closely in developing this project.

Authors' Contact Information: Andrea Botticella, andrea.botticella@studenti.polito.it, Politecnico di Torino, Turin, Italy; Elia Innocenti, elia.innocenti@studenti.polito.it, Politecnico di Torino, Turin, Italy; Renato Mignone, renato.mignone@studenti.polito.it, Politecnico di Torino, Turin, Italy; Simone Romano, simone.romano2@studenti.polito.it, Politecnico di Torino, Turin, Italy.

The attacks are relatively well-distributed throughout the day, but notable peaks occur between 19:00 and 6:00, with lower activity observed from 9:00 to 16:00. This pattern might be explained by the likelihood that attackers schedule their activities during non-working hours when individuals and organizations are less likely to monitor or respond to security incidents.

#### A.1.2 Attack Frequency by Month

This plot illustrates the distribution of SSH attacks across different months. It highlights seasonal trends, showing months with higher attack frequencies. Such insights can be useful for anticipating periods of increased security threats and allocating resources accordingly.

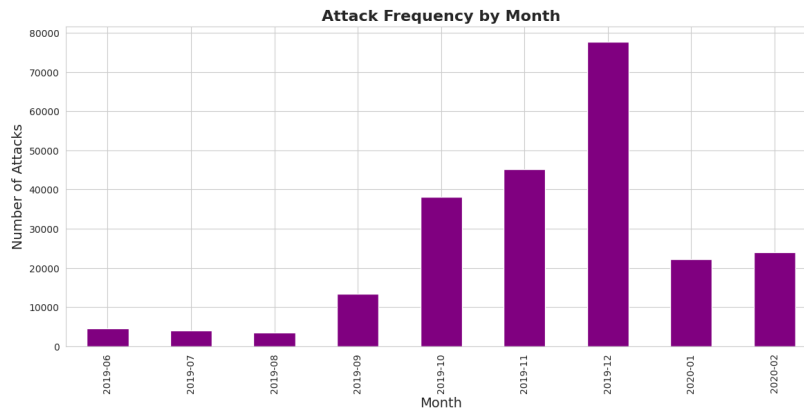


Fig. 2. Distribution of SSH attacks by month - The plot reveals seasonal trends in attack frequency

The number of attacks shows a noticeable increase from October to December, followed by a return to lower and more constant levels in the subsequent months. This trend raises the question of whether a seasonal pattern exists, possibly linked to factors such as end-of-year activities, holidays, or specific campaigns by attackers during this period.

#### A.1.3 Attack Frequency by Year

The plot compares the frequency of SSH attacks between the years 2019 and 2020.

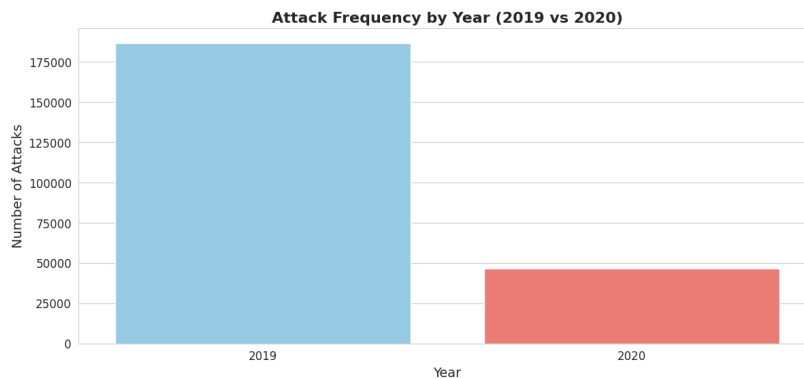


Fig. 3. Distribution of SSH attacks by year - The plot shows the overall trend of attacks over multiple years

#### A.1.4 Temporal Series of SSH Attacks

This time series plot illustrates the distribution of SSH attacks over time, providing a clear view of how attack frequency fluctuates throughout the dataset. When combined with other plots, such as those highlighting daily or seasonal patterns, it enables a more comprehensive temporal analysis, helping to uncover deeper insights into attack trends and dynamics.

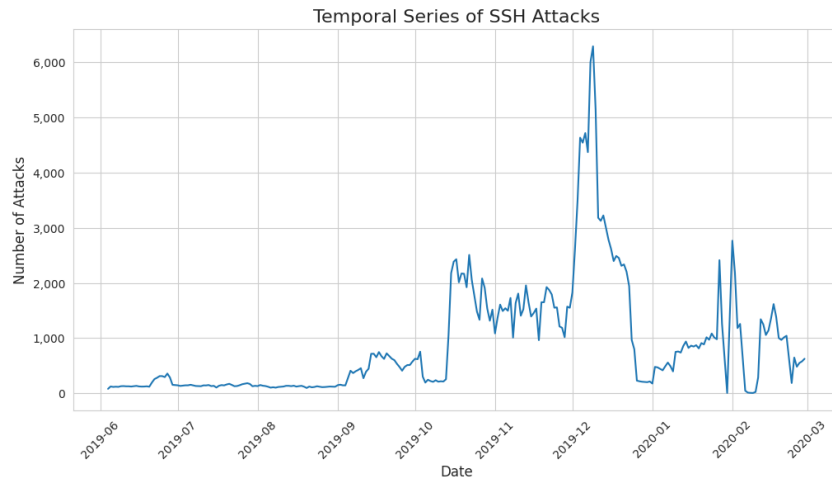


Fig. 4. Time series plot of SSH attacks over the entire dataset - The plot provides a view of attack patterns over time

#### A.1.5 Intents Over Timestamps

The plot visualizes the distribution of different attack intents over time. It categorizes the intents: Defense Evasion, Harmless, Impact, Discovery, Persistence, Execution, and Others. By analyzing this plot, we can identify which intents are more prevalent at specific times, helping to prioritize security measures based on the nature of the threats.

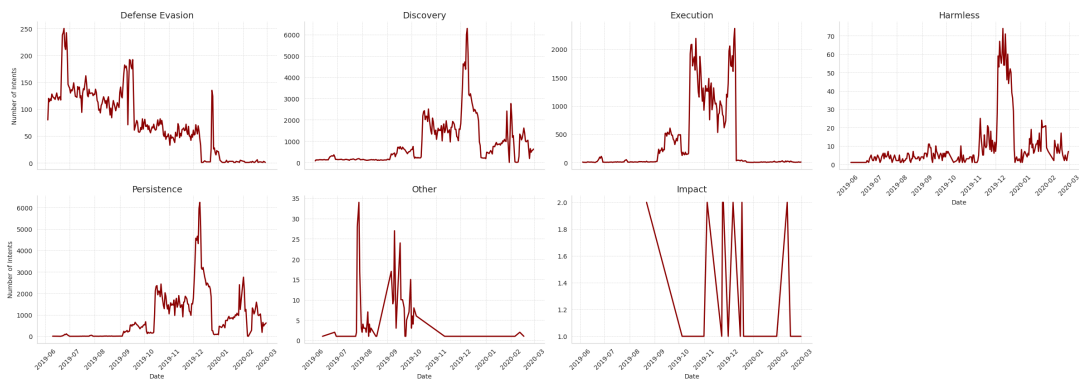


Fig. 5. Attack intents over timestamps - The plot provides insights into the temporal patterns of different attack intents

## B SUPERVISED LEARNING - CLASSIFICATION

### B.1 Logistic Regression

In this section, we detail the steps and results obtained from using the Logistic Regression model during the supervised learning phase of the project. Logistic Regression served as a baseline model to provide an initial understanding of the classification problem. Despite its simplicity, it offered valuable insights into the multi-label classification task.

#### B.1.1 Model Training

The Logistic Regression model was trained using its default configuration. Specifically, the `lbfgs` solver was utilized with a regularization parameter  $C = 1$ . The training process aimed to identify potential overfitting or underfitting issues and establish baseline performance metrics.

The dataset was preprocessed using the TF-IDF representation of the session texts, which assigned weights to words based on their frequency and relevance within the dataset. Multi-label binary encoding was applied to the `Set_Fingerprint` column to ensure compatibility with the model.

#### B.1.2 Evaluation Metrics

The Logistic Regression model was evaluated using standard classification metrics, including weighted F1-scores, precision, and recall. The evaluation metrics highlighted the strengths and weaknesses of the model in handling imbalanced classes.

PERFORMANCE ON TRAIN SET: Logistic Regression					
Model	Set	Attack	Precision	Recall	F1-Score Accuracy
Logistic Regression Train	Defense Evasion		0.993876	0.983690	0.988718 0.996641
Logistic Regression Train	Discovery		0.947280	0.921436	0.933970 0.999001
Logistic Regression Train	Execution		0.995354	0.993764	0.994543 0.994771
Logistic Regression Train	Harmless		0.985379	0.946215	0.962221 0.991436
Logistic Regression Train	Impact		0.499923	0.500000	0.499962 0.999847
Logistic Regression Train	Other		0.997616	0.973211	0.985106 0.999920
Logistic Regression Train	Persistence		0.998091	0.992316	0.995183 0.998375

PERFORMANCE ON TEST SET: Logistic Regression					
Model	Set	Attack	Precision	Recall	F1-Score Accuracy
Logistic Regression Test	Defense Evasion		0.993304	0.981371	0.987249 0.996267
Logistic Regression Test	Discovery		0.956802	0.922475	0.938964 0.999156
Logistic Regression Test	Execution		0.994998	0.993154	0.994056 0.994321
Logistic Regression Test	Harmless		0.995662	0.947761	0.965015 0.991332
Logistic Regression Test	Impact		0.499986	0.500000	0.499993 0.999971
Logistic Regression Test	Other		0.994971	0.980575	0.987667 0.999928
Logistic Regression Test	Persistence		0.998024	0.992011	0.994995 0.998326

Fig. 6. Evaluation Metrics for Logistic Regression Model

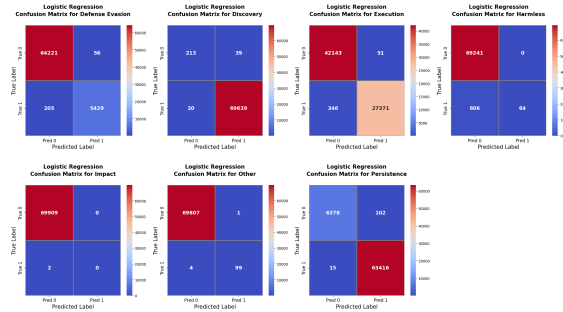


Fig. 7. Confusion Matrix for Logistic Regression Model

The confusion matrix provided a breakdown of true positives, false positives, false negatives, and true negatives for each intent. Figure 11 shows the confusion matrix for the Logistic Regression model.

#### B.1.3 Hyperparameter Tuning

Grid search was performed to optimize the Logistic Regression model's hyperparameters. The search focused on varying the regularization parameter  $C$  over a range of values  $[0.1, 1, 10, 100]$  to identify the configuration that maximized weighted F1-scores.

The optimized model exhibited improved performance compared to the baseline, particularly for intents with smaller sample sizes. Figure 10 illustrates the weighted F1-scores for different values of  $C$ .

#### B.1.4 Comparative Analysis of Baseline and Optimized Models

The optimized Logistic Regression model demonstrated a moderate improvement in precision and recall compared to the baseline. However, its overall performance remained slightly inferior to more complex models like Random

Forest and SVM. The comparative analysis underscores the importance of selecting models suited to the dataset's characteristics and problem requirements.

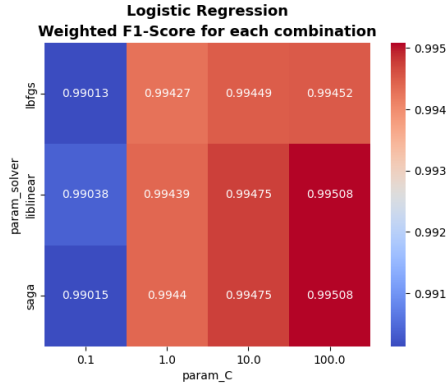


Fig. 8. Weighted F1-Scores for Hyperparameter Tuning

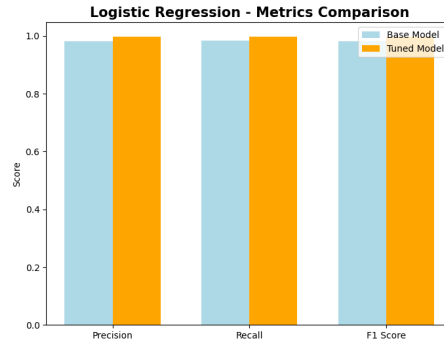


Fig. 9. LR Model Comparison

## B.2 Random Forest and SVM

### B.2.1 Comparative Analysis of Baseline and Optimized Models

Here we compare the baseline and optimized Random Forest and Support Vector Machine (SVM) models using metrics such as accuracy, precision, recall, F1-score, and AUC-ROC to evaluate improvements from hyperparameter optimization.

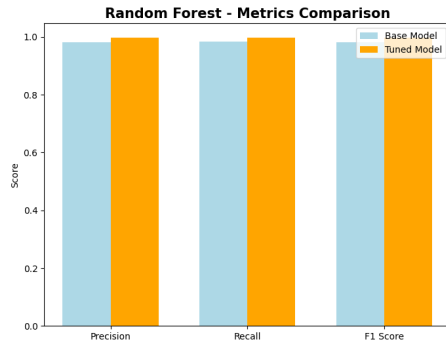


Fig. 10. RF Model Comparison



Fig. 11. SVM Model Comparison

### B.2.2 Probability Distributions of Predictions

Analyzing the probability distributions of predictions provides insights into the confidence levels of each model. Figures 12 and 13 depict the probability distributions for the Random Forest and SVM models, respectively. These visualizations help identify cases where the models are uncertain, which is critical for understanding their behavior.

The Random Forest model shows a bimodal distribution, with probabilities clustering near 0 and 1. This indicates high confidence for most predictions, with fewer uncertain cases. However, a small subset of predictions near 0.5 suggests some ambiguity in classification.

The SVM model, on the other hand, displays a more spread-out probability distribution, reflecting its sensitivity to margin violations and reliance on support vectors. This behavior could lead to better generalization but might introduce uncertainty in certain predictions.

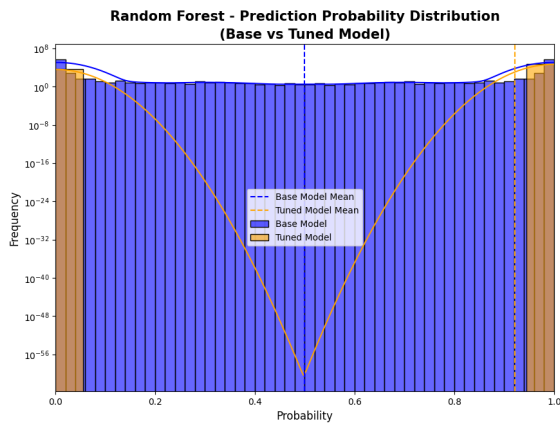


Fig. 12. Random Forest Probability Distribution

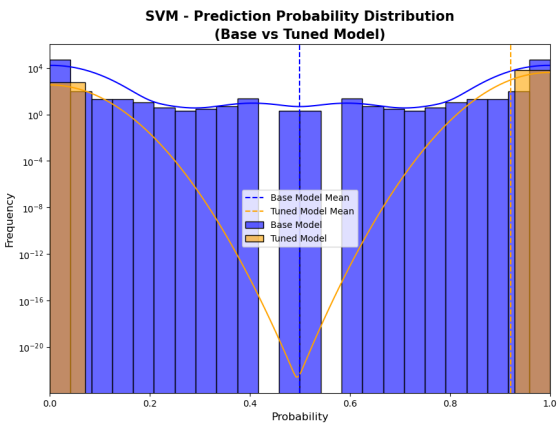


Fig. 13. SVM Probability Distribution

C UNSUPERVISED LEARNING - CLUSTERING

C.1 text  
text

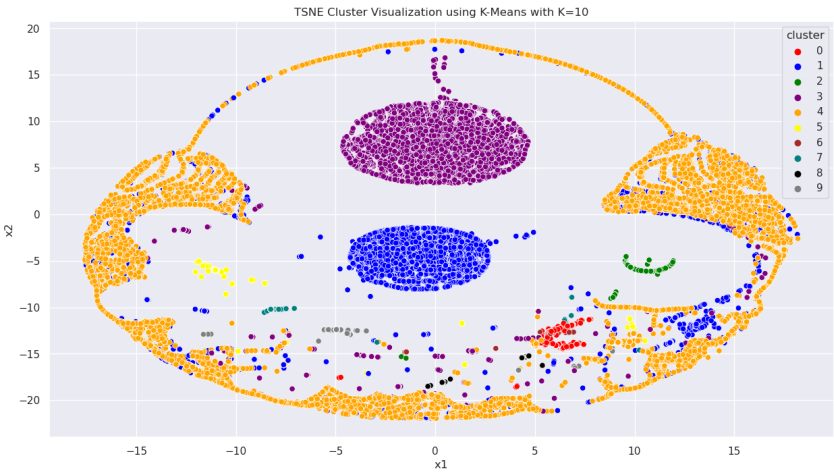


Fig. 14. t-SNE Visualization of K-Means Clusters

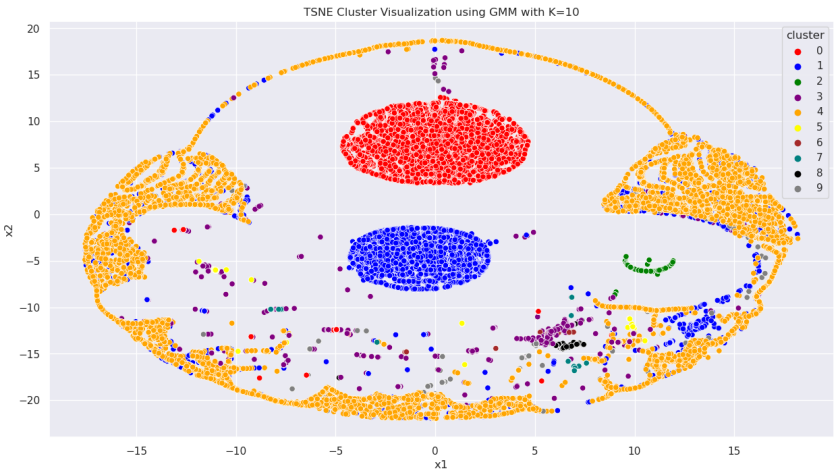


Fig. 15. t-SNE Visualization of GMM Clusters

C.2 matrices  
text

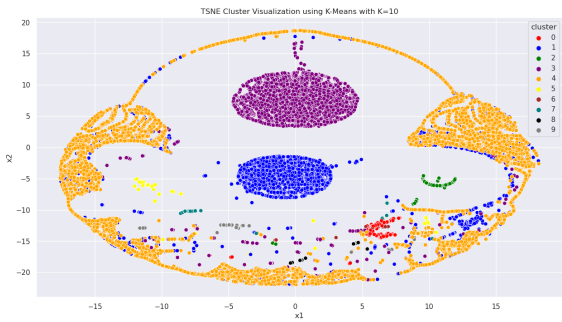


Fig. 16

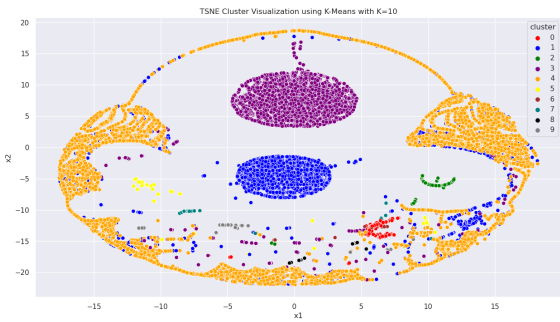


Fig. 17



## D LANGUAGE MODEL EXPLORATION

### D.1 Training Configuration

The model was implemented using BERT (bert-base-uncased) with the following key configurations:

- Maximum sequence length: 128 tokens
- Learning rate:  $4e-5$  with AdamW optimizer ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$ ,  $\epsilon = 1e-6$ )
- Training epochs: 4
- Gradient accumulation steps: 4
- Mixed precision training enabled
- Linear learning rate scheduler without warmup
- Loss function: Binary Cross-Entropy with Logits

### D.2 Model Performance Metrics

*D.2.1 Class-wise F1 Scores* The model demonstrates varying performance across classes (Figure 18):

- Excellent performance ( $F1 \geq 0.98$ ) for Defense Evasion, Discovery, Execution, Other, and Persistence classes
- Perfect scores ( $F1 = 1.00$ ) for Discovery and Persistence
- Significantly lower performance for Harmless class ( $F1 = 0.22$ )
- Impact class shows minimal detection capability ( $F1 = 0.00$ )

*D.2.2 Detailed Performance Metrics* The performance metrics (Figure 19) reveal:

- Most classes achieve balanced precision and recall scores
- The Harmless class shows a significant disparity between precision and recall
- The Impact class shows minimal performance across all metrics

*D.2.3 Precision-Recall Analysis* The Precision-Recall curves (Figure 20) demonstrate:

- Most classes maintain high precision ( $>0.95$ ) across different recall thresholds
- The "Other" category shows a sharp decline in precision at approximately 0.2 recall
- Impact and Harmless classes demonstrate poor precision-recall trade-offs
- Defense Evasion, Discovery, Execution, and Persistence maintain near-perfect precision until very high recall values

*D.2.4 Prediction Probability Distribution* The prediction probability histograms (Figure 21) reveal:

- Most classes exhibit a strong binary separation in prediction probabilities
- Defense Evasion, Discovery, and Execution show high-confidence predictions clustered near 0 and 1
- The Persistence class shows a similar pattern but with a smaller proportion of low-probability predictions
- Harmless and Impact classes show predominantly low-probability predictions, indicating potential class imbalance issues

### D.3 Recommendations for Model Improvement

Based on the analysis, several potential improvements could be considered:

*D.3.1 Class Imbalance Mitigation*

- Implement class weights or sampling techniques for Harmless and Impact classes
- Consider data augmentation for underrepresented classes

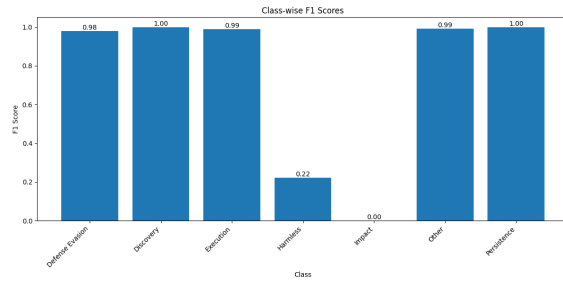


Fig. 18. F1 scores across different classes showing the model's classification performance for each category.

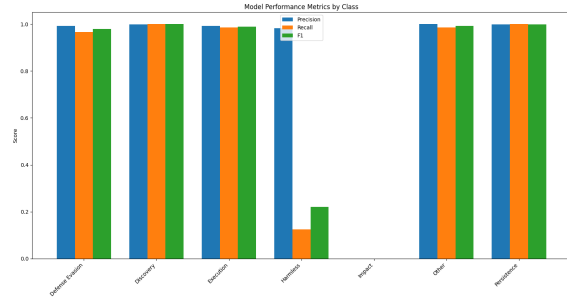


Fig. 19. Detailed breakdown of Precision, Recall, and F1 scores for each class.

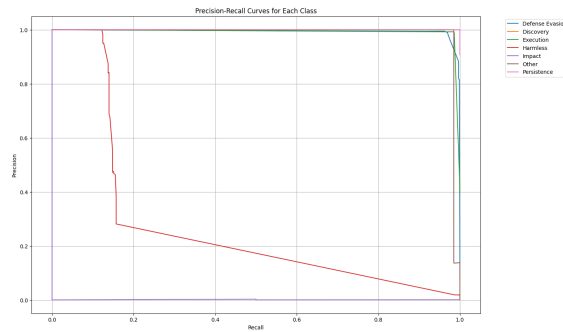


Fig. 20. Precision-Recall curves for each class showing the trade-off between precision and recall at different classification thresholds.

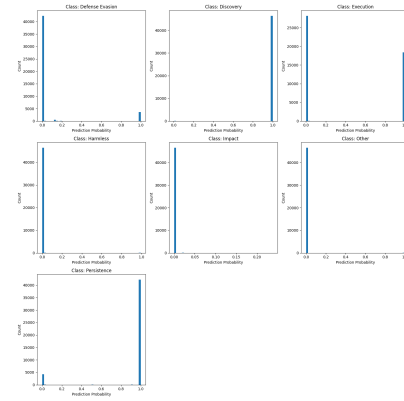


Fig. 21. Distribution of prediction probabilities for each class showing the model's confidence in its predictions.

### D.3.2 Model Architecture

- Experiment with different BERT variants
- Consider ensemble approaches for improving performance on challenging classes

### D.3.3 Training Strategy

- Implement curriculum learning for difficult classes
- Explore different learning rate schedules
- Consider longer training with early stopping

### D.3.4 Data Quality

- Review and potentially relabel samples in the Impact class
- Analyze misclassified examples in the Harmless class