

Laboratory 1 Report

ANDREA BOTTICELLA*, s347291, Politecnico di Torino, Italy

ELIA INNOCENTI*, s345388, Politecnico di Torino, Italy

SIMONE ROMANO*, s344024, Politecnico di Torino, Italy

This work presents the development of Feed-Forward Neural Networks (FFNNs) for intrusion detection using the CICIDS2017 dataset. Through six progressive tasks, we address data cleaning, feature bias, class imbalance, architectural design, and overfitting control. The final model—a three-layer FFNN (128-64-32) with ReLU activation, AdamW optimization, and Batch Normalization—achieved 96% accuracy and 0.94 macro-F1. Results highlight that bias mitigation, weighted loss optimization, and regularization are essential for reliable AI-based IDS capable of detecting diverse network attacks in real-world cybersecurity contexts.

1 INTRODUCTION

This laboratory explores the implementation of a **Feed-Forward Neural Network (FFNN)** using the **CICIDS2017** dataset, a standard benchmark for intrusion detection research. The goal is to construct a complete machine learning pipeline in PyTorch—from raw data preparation to model evaluation—to analyze how architectural choices and preprocessing strategies affect classification performance in cybersecurity contexts.

The experiment unfolds through six progressive tasks that build complexity systematically. Beginning with **data preprocessing**—encompassing cleaning, scaling, and outlier management—the work advances to **baseline FFNN training** using a single hidden layer architecture. Subsequently, **feature bias analysis** examines the influence of specific attributes like *Destination Port*, followed by **loss-function weighting** to address class imbalance challenges. The investigation then explores **deep network optimization** through architectural and optimizer variations, culminating in **regularization** strategies including dropout, batch normalization, and weight decay techniques.

2 TASK 1 — DATA ANALYSIS AND PREPROCESSING

2.1 Dataset Overview

The raw dataset contained **31,507 samples** and **17 features**, including numerical flow statistics and a categorical label identifying traffic types: *Benign* ($\approx 63\%$ of the samples), *DoS Hulk*, *PortScan*, and *Brute Force*. A class distribution plot (Figure 1) visually confirmed this imbalance, motivating the use of **stratified splitting** later in the pipeline.

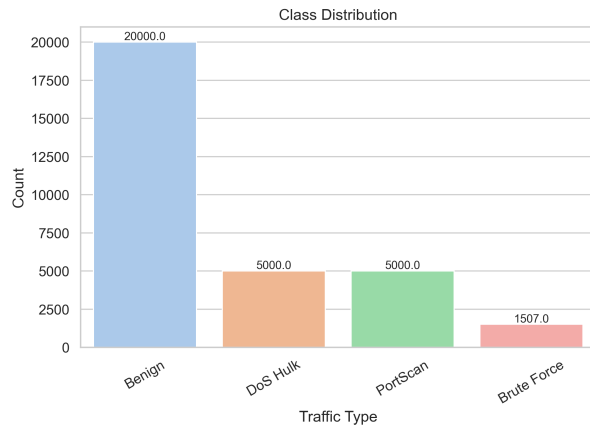


Fig. 1. Class distribution in the raw dataset.

Step	Removed	Remaining
Drop NaN	20	31,487
Drop duplicates	2,114	29,393
Drop infinite values	7	29,386

Table 1. Summary of the raw dataset.

Before any training, the dataset must be cleaned and normalized to ensure the model learns meaningful statistical relationships rather than artifacts of data noise or imbalance. The preprocessing pipeline addresses three critical aspects: **data quality** through the removal of missing, duplicate, and infinite values; **class imbalance** by preserving

*The authors collaborated closely in developing this project.

proportional representation during splitting; and **feature scaling** to standardize feature ranges and stabilize neural training convergence.

2.2 Data Cleaning and Partitioning

Table 1 summarizes the cleaning process: duplicate rows, missing values (NaN), and infinite entries were systematically removed. Thus, the dataset was reduced to **29,386 samples**, meaning **2,121 rows** were discarded (2,114 between missing and duplicates, and 7 infinite values). Categorical labels were then encoded numerically: *Benign* = 0, *Brute Force* = 1, *DoS Hulk* = 2, *PortScan* = 3.

The cleaned data were divided into training (60%, 17,631 samples), validation (20%, 5,877 samples), and test (20%, 5,878 samples) subsets using **stratified sampling** with a fixed seed, ensuring consistent class proportions and reproducibility. Exploration of numerical attributes revealed high variability and heavy-tailed distributions (Figures 2), indicating the presence of outliers and the need for normalization.

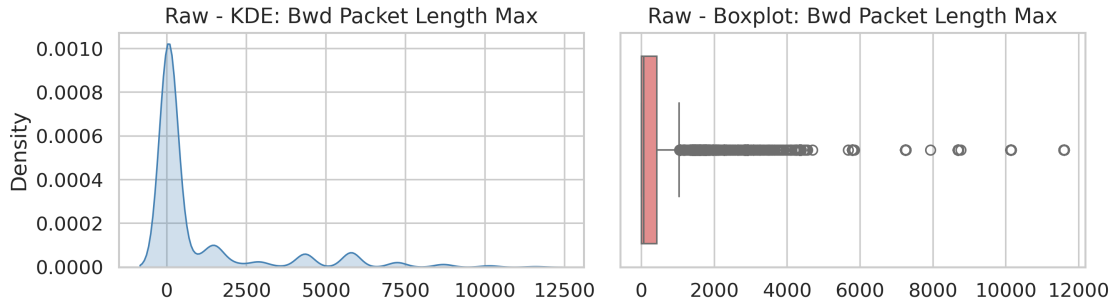


Fig. 2. Examples of feature (Bwd Packet Length Max) kernel density plot and boxplot, highlighting outliers.

2.3 Outlier Detection and Normalization

Outliers were analyzed through the **Z-score** and **IQR** methods. Both confirmed extreme values in features such as *Bwd Packet Length Max*, *Flow Duration*, and *Fwd IAT Std*. Outliers were retained to preserve data realism, and normalization was applied to reduce their effect.

Two scaling methods were tested: **StandardScaler** and **RobustScaler**. The density comparison plots (Figure 3) showed that both methods effectively normalized distributions, but **RobustScaler** produced more compact, less skewed curves for highly variable features. However, the numerical statistics and subsequent training experiments revealed negligible performance difference between the two. Therefore, **StandardScaler** was ultimately adopted for simplicity and interpretability, offering smoother loss curves in preliminary trials.

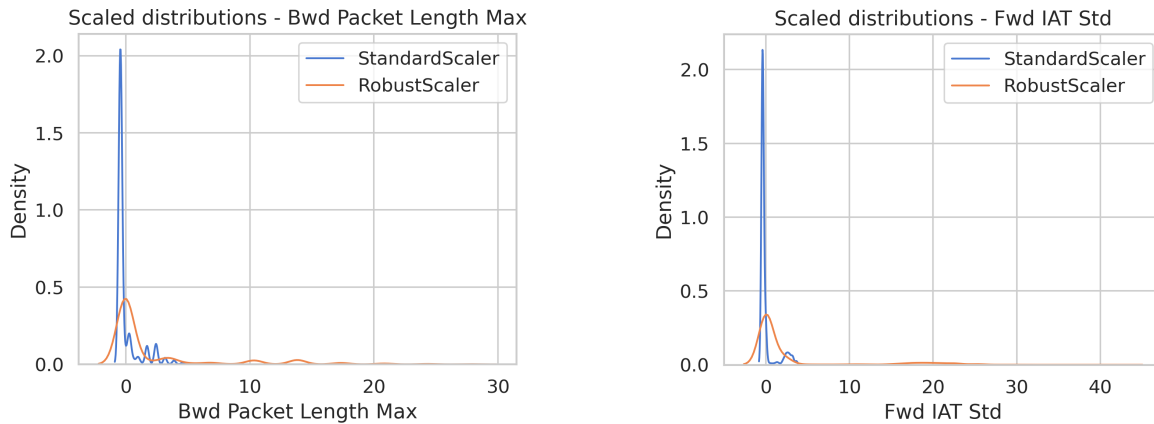


Fig. 3. Comparison of normalization effects using StandardScaler and RobustScaler.

This preprocessing ensured consistent, balanced, and properly scaled data, forming a robust foundation for training the Feed-Forward Neural Network.

3 TASK 2 — SHALLOW FEED-FORWARD NEURAL NETWORK (FFNN)

3.1 Model Configuration

A **single-layer Feed-Forward Neural Network (FFNN)** was trained with 32, 64, and 128 neurons to study the effect of network size on learning and generalization. Each model used the Adam optimizer ($\text{lr} = 0.0005$), linear activation function, and early stopping over 100 epochs, minimizing categorical cross-entropy on the same partitions defined in Task 2. Inputs are standardized features, while outputs are logits over the four classes.

3.2 Training Dynamics

The training and validation loss curves (Figures ??–??) show consistent convergence for all models, with rapid loss reduction during early epochs followed by stable plateaus. No overfitting was observed. All models converged to similar validation loss levels (~ 0.29), with the 64-neuron model early stopping at 100 epochs.

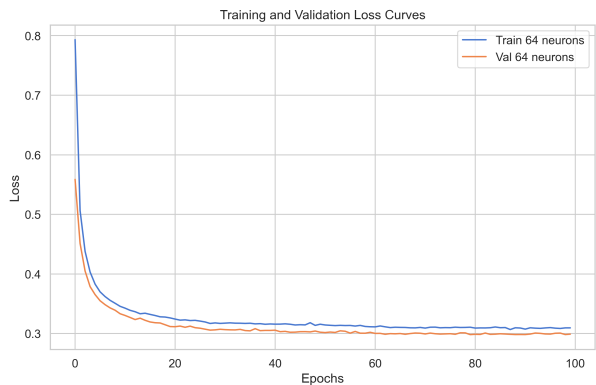


Fig. 4. 64 neurons

Class	Precision	Recall	F1-Score
0 Benign	0.8872	0.9524	0.9187
1 Brute Force	0.0000	0.0000	0.0000
2 DoS Hulk	0.9854	0.8758	0.9274
3 Port Scan	0.8270	0.8918	0.8581

Table 2. Full validation report for the 64-neuron model.

3.3 Validation Performance Analysis

Model	Accuracy	F1 (macro avg)	F1 (weighted avg)
32 Neurons	0.8860	0.6761	0.8651
64 Neurons	0.8860	0.6760	0.8651
128 Neurons	0.8802	0.6709	0.8601

Table 3. Validation metrics for linear models.

Besides loss trajectories, validation metrics (Table 3) highlight performance differences among the three configurations. The **32** and **128** neuron models achieved good overall accuracy (~ 0.88) but failed on the minority *Brute Force* class (precision and recall = 0), primarily learning majority classes (*Benign* and *Port Scan*). The **64**-neuron model matched this accuracy while slightly improving macro F1 (0.6760 vs. 0.6709 for 128 neurons) by better balancing class predictions. Considering both the loss trajectories and the validation metrics for this run, the **64**-neuron model was selected for detailed class-wise analysis.

3.4 Activation Function Study and Generalization

Replacing the linear activation with ReLU (64 neurons) accelerated convergence and markedly improved minority-class recognition. In particular, the *Brute Force* class (1) F1 rose from 0.28 (linear model; Table 8) to 0.77 with ReLU (Table 4), showing that ReLU helped capture more complex patterns for this rare attack type. Figure 5 illustrates the faster/stabler loss dynamics, and validation/test metrics remained closely aligned, confirming good generalization. Overall the model performs best on *Benign*, *DoS Hulk* and *Port Scan* ($\text{F1} = 0.96, 0.95, 0.92$ respectively) while now also handling *Brute Force* effectively, indicating a strong across-the-board improvement.

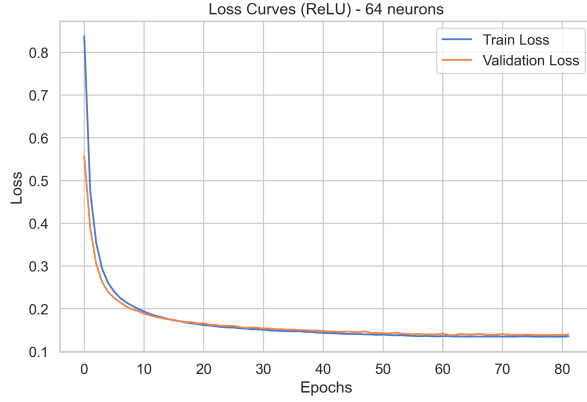


Fig. 5. Loss curves using ReLU activation.

Class	Precision	Recall	F1-Score
0 <i>Benign</i>	0.9593	0.9675	0.9634
1 <i>Brute Force</i>	0.7768	0.9371	0.8494
2 <i>DoS Hulk</i>	0.9986	0.9043	0.9491
3 <i>Port Scan</i>	0.9380	0.9196	0.9287
Accuracy	Macro F1	Weighted F1	
0.9498	0.9226	0.9502	

Table 4. Validation set report.

4 TASK 3 — IMPACT OF SPECIFIC FEATURES (DESTINATION PORT)

4.1 Dataset Bias and Feature Dependence

During the dataset inspection, it was observed that most of the **Brute Force** attacks shared the same **Destination Port** value (port 80). This is an unrealistic scenario, as *Brute Force* attacks can occur on any service requiring authentication. This systematic bias introduces a wrong inductive pattern, leading the model to associate port 80 exclusively with Brute Force traffic instead of learning meaningful behavioral features. As a result, the model risks overfitting to an artifact of data collection rather than generalizing to real-world cases.

4.2 Evaluating Bias via Port Substitution

To quantify this effect, the trained 64-neuron ReLU model was evaluated on a **modified test set** in which all *Brute Force* samples had their destination port changed from 80 to 8080. While the model performed well on the original validation set (*Brute Force*: F1 ~0.85, overall accuracy ~95%), its performance degraded sharply on the modified test set (*Brute Force*: F1 ~0.08, overall accuracy ~90%). This dramatic decline confirms that the model learned a spurious dependency on the port feature, failing to recognize attacks when this shortcut was removed.

4.3 Effect of Removing the Destination Port Feature

To mitigate this bias, the destination port attribute was excluded, and the dataset was reprocessed. After duplicate and NaN removal, the number of *PortScan* samples **decreased drastically** from 5,000 to only 285, as shown in Figure 6. This reduction indicates that many *PortScan* flows were nearly identical except for their port values; removing the feature exposed these redundancies.



Fig. 6. Updated class distribution.

Class	Precision	Recall	F1-Score
0 <i>Benign</i>	0.8997	0.9667	0.9320
1 <i>Brute Force</i>	0.1630	0.0526	0.0796
2 <i>DoS Hulk</i>	0.9971	0.8928	0.9421
3 <i>Port Scan</i>	0.9300	0.9175	0.9237
Accuracy	Macro F1	Weighted F1	
0.9046	0.7193	0.8906	

Table 5. Test set report - port changed before scaling.

4.4 Class Balance Considerations

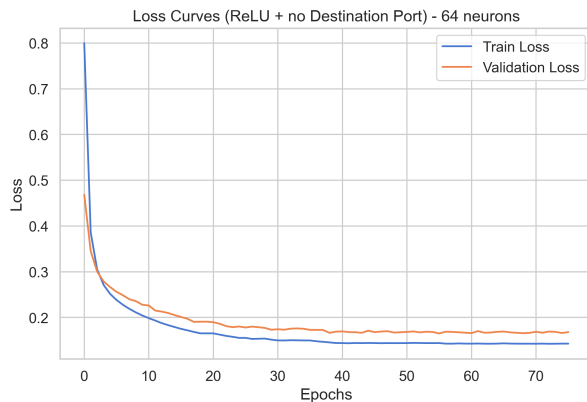
Even after preprocessing, the dataset remains **imbalanced**, with benign samples far exceeding attack ones and minority classes (*Brute Force*, *Port Scan*) underrepresented. Although dropping the destination port feature improves robustness against spurious correlations, addressing class imbalance remains necessary to prevent the model from favoring majority classes.

In summary, this task highlights how biased or overly discriminative features can mislead model learning. Consequently, removing or down-weighting the Destination Port attribute is justified for the subsequent experiments. Its exclusion produces a more reliable though smaller dataset, forcing the model to learn from intrinsic traffic behaviors rather than arbitrary identifiers.

5 TASK 4 — THE IMPACT OF THE LOSS FUNCTION (CLASS WEIGHTED)

5.1 Removing the Destination Port Feature

The best-performing model from previous tasks (64 neurons, ReLU activation) was retrained after excluding the *Destination Port* feature to eliminate the bias discussed earlier. This modification had a mixed effect: overall accuracy remained stable, and performance on the *Brute Force* class slightly improved, confirming that the model no longer relied on a biased feature. However, the ability to recognize the rarest class, *PortScan*, declined significantly (F1-score dropped from 0.92 to 0.21), suggesting that the model had previously exploited this feature as a strong shortcut for *PortScan* detection.



Class	Precision	Recall	F1-Score
0 Benign	0.9455	0.9751	0.9601
1 Brute Force	0.7927	0.9091	0.8469
2 DoS Hulk	0.9880	0.8486	0.9130
3 Port Scan	0.4444	0.1404	0.2133
Accuracy	0.9386	Macro F1	0.7333
		Weighted F1	0.9353

Table 6. Test set report.

Fig. 7. Loss curves for the best model (no Destination Port).

5.2 Class Weights and Weighted Loss

To counter the imbalance identified in previous tasks, the strategy was to apply class weights in the loss function. Weighted cross-entropy was then adopted to penalize misclassifications of minority classes more strongly, encouraging the model to learn balanced decision boundaries.

Benign	0.3326	Brute Force	3.9372	DoS Hulk	1.4521	Port Scan	19.7091
---------------	--------	--------------------	--------	-----------------	--------	------------------	---------

Table 7. Class weights used for weighted cross-entropy loss.

These weights were estimated on the training partition to avoid data leakage. Estimating class weights from training data ensures that information from validation or test sets is not used during training or weight calculation, and allows the weighted loss to emphasize minority classes while preserving evaluation integrity.

5.3 Effect of Weighted Cross-Entropy

Training with the weighted loss produced smoother convergence and more balanced class performance (Figure 8). Overall accuracy decreased slightly (accuracy 0.9386 → 0.9243), while recall for underrepresented classes improved markedly — *PortScan* recall rose from 0.1404 to 0.8421 and *Brute Force* recall increased slightly from 0.9091 to 0.9545. This confirms that the weighted cross-entropy promotes fairness across classes by reducing bias toward dominant traffic types, at the cost of a small drop in global accuracy.

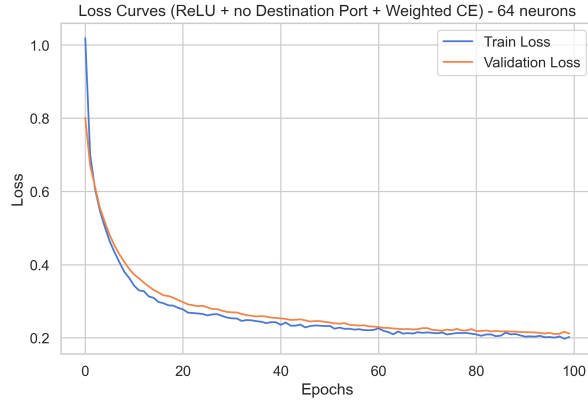


Fig. 8. Loss curves (weighted cross-entropy).

Class	Precision	Recall	F1-Score
0 <i>Benign</i>	0.9800	0.9275	0.9530
1 <i>Brute Force</i>	0.7398	0.9545	0.8336
2 <i>DoS Hulk</i>	0.9459	0.9056	0.9253
3 <i>Port Scan</i>	0.2553	0.8421	0.3918
Accuracy	Macro F1	Weighted F1	
0.9243	0.7759	0.9335	

Table 8. Test set report.

In summary, applying a class-weighted loss increased sensitivity to underrepresented attacks, but this came at the expense of precision (*PortScan* precision $0.4444 \rightarrow 0.2553$). This experiment highlights the trade-off between improving per-class recall for minority classes and reducing precision and global performance.

6 TASK 5 — DEEP NEURAL NETWORKS, BATCH SIZE, AND OPTIMIZERS

6.1 Architecture Depth Analysis

In this task, the Feed Forward Neural Network (FFNN) was extended to deeper configurations to evaluate the effect of architectural depth on classification performance. Six architectures ($[16,8,4]$, $[32,16,8]$, $[32,16,8,4]$, $[16,16,8,8]$, $[32,32,16,8,4]$, $[32,32,8,16,16]$) were trained with depths ranging from three to five hidden layers, with variable neuron widths per layer (2-32). All models used the ReLU activation, the AdamW optimizer ($\text{lr} = 5e-4$), batch size 64, and early stopping with patience = 20 and min_delta = $1e-5$.

Each model was trained for up to 50 epochs. All exhibited smooth and consistent convergence, as both training and validation losses decreased stably before plateauing, indicating successful optimization without overfitting (see Figure 9).

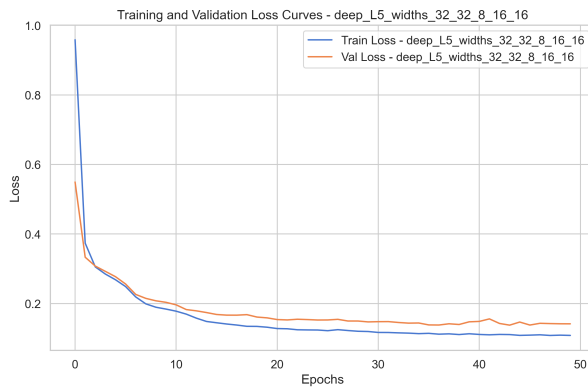


Fig. 9. Losses for the best-performing architecture.

Class	Precision	Recall	F1-Score
0 <i>Benign</i>	0.9664	0.9805	0.9734
1 <i>Brute Force</i>	0.8390	0.9476	0.8900
2 <i>DoS Hulk</i>	0.9858	0.8952	0.9383
3 <i>Port Scan</i>	0.8571	0.6316	0.7273
Accuracy	Macro F1	Weighted F1	
0.9593	0.8822	0.9589	

Table 9. Test set report.

The second architecture with 5 layers ($[32,32,8,16,16]$) achieved the best results on the validation set with an accuracy of **95.86%**, weighted F1 of **0.9581**, and macro F1 of **0.8623**. This model also stood out as one of the few models capable of effectively detecting the minority class (class 3, $F1 = 0.64$), while maintaining high performance on the dominant classes. On the test set, it achieved a **95.93%** accuracy and macro F1 of **0.8822**, confirming good generalization with improved minority-class recall (0.7273).

The results show that deeper architectures yield better feature abstraction and improved robustness, particularly for imbalanced datasets. The minority class improved significantly without degrading the performance on the majority classes, suggesting well-formed and stable decision boundaries.

6.2 Effect of Batch Size

To analyze the impact of batch size, the best-performing architecture was retrained with batch sizes {4, 64, 256, 1024}. Validation performance varied significantly with batch size, as shown in Table 11.

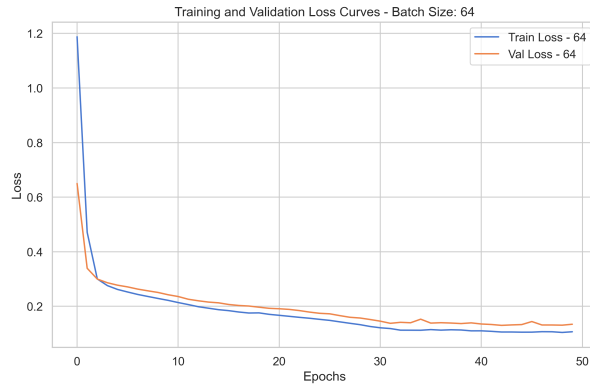


Fig. 10. Training and validation loss with batch size 64.

The optimal configuration was found at batch size 64 (see Table 10). Very small batches introduced excessive gradient noise, harming convergence, while very large batches led to overly smooth gradient estimates, biasing the model toward majority classes. Moderate batch sizes balanced noise and stability, allowing minority class learning. Training time decreased with larger batches due to fewer updates per epoch: 75.4 s (batch 4) vs. 3.1 s (batch 256).

Batch	Accuracy	Macro F1	Class 3 F1
4	0.9021	0.4864	0.0000
64	0.9591	0.8040	0.3689
256	0.9490	0.6888	0.0000
1024	0.8970	0.4629	0.0000

Table 11. Validation results for different batch sizes.

Class	Precision	Recall	F1-Score
0 <i>Benign</i>	0.9633	0.9864	0.9747
1 <i>Brute Force</i>	0.9276	0.9439	0.9357
2 <i>DoS Hulk</i>	0.9871	0.8915	0.9369
3 <i>Port Scan</i>	0.4130	0.3333	0.3689

Accuracy	Macro F1	Weighted F1
0.9591	0.8040	0.9580

Table 10. Validation set report.

Optimizer	Accuracy	Macro F1	Class 3 F1
SGD	0.7517	0.2146	0.0000
SGD (0.1)	0.7517	0.2146	0.0000
SGD (0.5)	0.7517	0.2146	0.0000
SGD (0.9)	0.8976	0.4634	0.0000
AdamW	0.9537	0.7852	0.3125

Table 12. Validation results for different optimizers.

6.3 Effect of Optimizer

The optimizers compared were SGD, SGD with momentum (0.1, 0.5, 0.9), and AdamW. Results confirmed that optimizer choice strongly affects both convergence rate and class balance (Figure 11).

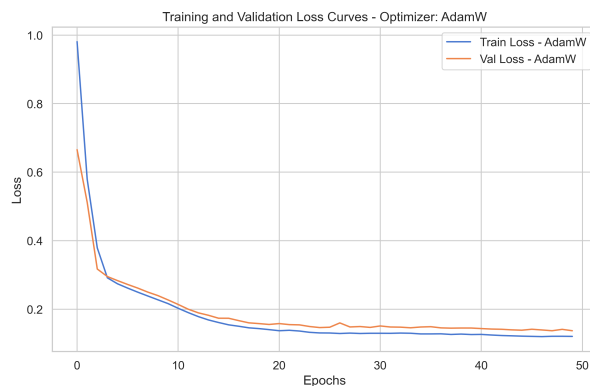


Fig. 11. Loss curve of the AdamW optimizer.

Class	Precision	Recall	F1-Score
0 <i>Benign</i>	0.9564	0.9867	0.9713
1 <i>Brute Force</i>	0.9308	0.9439	0.9373
2 <i>DoS Hulk</i>	0.9824	0.8643	0.9196
3 <i>Port Scan</i>	0.3846	0.2632	0.3125

Accuracy	Macro F1	Weighted F1
0.9537	0.7852	0.9519

Table 13. Validation set report.

SGD and its momentum variants often collapsed to majority-class predictions, showing low macro F1 (~0.21-0.46). In contrast, AdamW achieved **accuracy 0.9537**, **macro F1 0.7852**, and detected all classes.

It also provided the fastest convergence with minimal oscillation. Although slightly slower in wall-clock time (6.8 s vs. 5.3 s for SGD), AdamW produced the most balanced and generalizable model.

6.4 Learning Rate and Epochs Exploration

With AdamW and batch size 64 fixed, several learning rates and epoch limits were tested.

Very small learning rates ($\leq 1 \times 10^{-4}$) converge but the models collapse to majority-class predictions (no class 3 detections), even with more epochs. A moderate learning rate (5×10^{-4}) yields strong overall accuracy but still ignores the rarest class and can drift with longer training. A higher learning rate (5×10^{-3}) combined with early stopping enables minority-class learning without harming dominant-class performance.

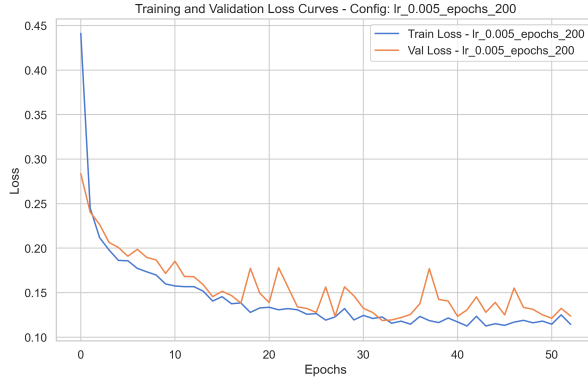


Fig. 12. Losses for the best model ($lr=0.005$, 200 epochs).

Class	Precision	Recall	F1-Score
0 <i>Benign</i>	0.9664	0.9873	0.9767
1 <i>Brute Force</i>	0.9281	0.9476	0.9377
2 <i>DoS Hulk</i>	0.9927	0.8797	0.9328
3 <i>Port Scan</i>	0.5000	0.5789	0.5366

Accuracy	Macro F1	Weighted F1
0.9611	0.8460	0.9611

Table 14. Test set report.

The best configuration used $lr = 0.005$ and $max_epochs = 200$, with early stopping at epoch 53. This setup improved minority-class recognition and achieved the highest macro F1 and overall accuracy. The corresponding loss curve is shown in Figure 12.

6.5 Summary of Findings

The deep architecture with 5 layers ([32,32,8,16,16]), AdamW optimizer, batch size 64, and learning rate 0.005 achieved the most balanced and generalizable results. Depth improved expressiveness; AdamW enabled stable convergence and robust minority detection. Proper batch sizing and learning rate tuning were crucial in mitigating the effects of class imbalance and preventing underfitting or overfitting.

7 TASK 6 — OVERFITTING AND REGULARIZATION TECHNIQUES

In this task, various regularization techniques were explored to mitigate overfitting, starting from a new baseline deep model with 6 layers ([256,128,64,32,16,8]), ReLU activations and the AdamW optimizer.

7.1 Baseline Model

The baseline deep model (AdamW, no explicit regularization) shows consistent convergence with both training and validation losses stabilizing (final Train Loss ≈ 0.1022 , Val Loss ≈ 0.1188 ; see Figure 14). The validation loss remains slightly higher than the training loss, indicating good generalization rather than strong overfitting. Both curves plateau together. The final validation and test accuracies are high (Validation accuracy = 96.24%, Test accuracy = 96.46%).

7.2 Effect of Normalization and Regularization

Several regularization strategies were applied to investigate their impact on convergence and performance:

- **Dropout (0.5):** Increased generalization pressure but produced weaker overall accuracy (Validation accuracy = 94.39%, Test accuracy = 93.93%). Final losses (Train ≈ 0.1462 , Val ≈ 0.1339) show the validation loss can be lower than the training loss; the minority class (class 3) was not predicted (zero precision/recall / F1 = 0.0).
- **Batch Normalization:** Produced less stable validation behaviour (final Train ≈ 0.1269 , Val ≈ 0.2123) with a higher validation loss at the end of training, suggesting sensitivity to batch statistics on this tabular dataset.

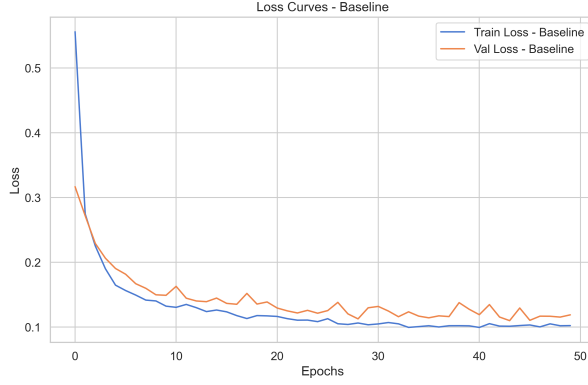


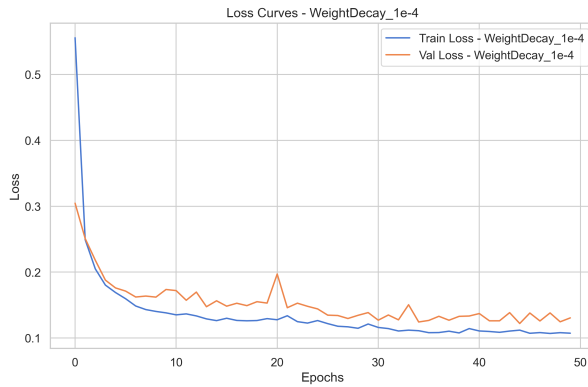
Fig. 13. Loss curves for the baseline model.

Class	Precision	Recall	F1-Score
0 <i>Benign</i>	0.9667	0.9876	0.9770
1 <i>Brute Force</i>	0.9249	0.9476	0.9361
2 <i>DoS Hulk</i>	0.9817	0.9017	0.9400
3 <i>Port Scan</i>	0.7750	0.5439	0.6392
Accuracy	Macro F1	Weighted F1	
0.9646	0.8731	0.9638	

Table 15. Test set report.

Validation/Test accuracies remained relatively high (Val = 95.51%, Test = 95.73%) but per-class recall for the minority class was degraded.

- **BatchNorm + Dropout (0.5)**: The combination appears to over-regularize the model: training and validation losses are higher than the baseline (final Train ≈ 0.1922 , Val ≈ 0.1737) and overall accuracy falls to $\approx 93\%$. Moreover, the model made no predictions for the minority class (F1 = 0.0), which suggests BatchNorm+Dropout harmed the decision boundary for rare classes (likely due to noisy batch statistics and excessive noise injection).
- **Weight Decay ($1e-4$)**: L2 regularization via weight decay yielded stable behaviour (final Train ≈ 0.1073 , Val ≈ 0.1305) and a good balance between stability and minority-class performance. Validation/Test accuracies remained close to the baseline (Val = 95.68%, Test = 95.86%) and recall for the minority class improved compared to all the other configurations (Baseline excluded).
- **Weight Decay + BN + Dropout (0.5)**: Combination of strong regularizers led to underfitting (final Train ≈ 0.1987 , Val ≈ 0.1708) and reduced accuracy (Val = 93.99%, Test = 93.55%); the minority class was again not detected in most runs.

Fig. 14. Loss curves using Weight Decay ($1e-4$).

Class	Precision	Recall	F1-Score
0 <i>Benign</i>	0.9695	0.9787	0.9741
1 <i>Brute Force</i>	0.9340	0.9439	0.9389
2 <i>DoS Hulk</i>	0.9588	0.9031	0.9301
3 <i>Port Scan</i>	0.3881	0.4561	0.4194
Accuracy	Macro F1	Weighted F1	
0.9568	0.8156	0.9572	

Table 16. Validation set report.

7.3 Final Observations

The experiments show a clear trade-off between overall performance (accuracy / weighted metrics) and robustness on the rare class.

The baseline produced the highest overall scores in these runs (Validation = 96.24%, Test = 96.46%) and also the best minority-class F1 among the configurations tested. **AdamW + light weight decay (1×10^{-4})** yielded slightly lower overall accuracy but improved stability and in some cases the minority-class recall compared to aggressive normalization schemes. Stronger regularizers (Dropout 0.5, BatchNorm combined with Dropout) tended to underfit this tabular task: they reduced overall accuracy and frequently led to the minority class being ignored (zero precision/recall in several runs).

If the project's primary objective is overall accuracy and weighted F1, the **baseline** is the preferred choice. However, if improving minority-class detection is critical, a mild regularization strategy like **Weight Decay (1e-4)** may be beneficial, despite a small drop in overall accuracy.

8 CONCLUSIONS

This laboratory comprehensively explored the design, optimization, and evaluation of Feed-Forward Neural Networks (FFNNs) for intrusion detection using the CICIDS2017 dataset. The six progressive tasks provided a systematic framework to address key challenges in data preprocessing, feature bias mitigation, class imbalance, model architecture, and regularization. The following summarizes the main findings from each stage and their broader implications.

- **Task 1 – Data Analysis and Preprocessing**

The initial task established a robust data foundation. Cleaning removed 2,121 redundant or invalid entries, ensuring data integrity. Exploratory analysis confirmed strong class imbalance and the presence of heavy-tailed distributions. Normalization via **StandardScaler** stabilized model training while preserving interpretability. This stage emphasized that well-structured, standardized input data is essential for stable convergence and reliable evaluation.

- **Task 2 – Shallow Feed-Forward Neural Network**

The single-layer FFNN served as a baseline for architectural exploration. Among hidden sizes of {32, 64, 128}, the 64-neuron configuration with ReLU activation achieved the best compromise between complexity and performance, reaching $\approx 95\%$ accuracy and a macro F1 of 0.92. The activation analysis revealed that ReLU substantially enhanced minority-class learning, correcting the failure observed with linear activation. This demonstrated how non-linear transformations are crucial for capturing subtle attack patterns.

- **Task 3 – Impact of Specific Features (Destination Port)**

This task exposed a critical dataset bias: all Brute Force samples shared the same destination port (80). When this feature was perturbed (changed to 8080), the model's Brute Force F1 dropped from 0.85 to 0.08, confirming an overreliance on non-generalizable cues. Removing the feature corrected this inductive bias but also revealed redundancy among PortScan samples, reducing their count from 5,000 to 285. This analysis underscored the importance of feature independence and dataset diversity to avoid misleading correlations that undermine generalization.

- **Task 4 – Impact of the Loss Function (Class Weighting)**

After removing the biased feature, class imbalance became more pronounced. Implementing a class-weighted cross-entropy loss improved recall for minority classes—especially PortScan (recall 0.14 \rightarrow 0.84)—with a limited accuracy drop (0.94 \rightarrow 0.92). This experiment highlighted the trade-off between global accuracy and fairness across classes, demonstrating that weighting the loss function is an effective yet delicate approach to improving attack diversity detection.

- **Task 5 – Deep Networks, Batch Size, and Optimizers**

Deeper architectures enhanced representation power and stability. The five-layer model ([32, 32, 8, 16, 16]) achieved the best results: accuracy 95.9%, macro F1 0.88, and improved minority-class F1 (0.73). Batch size analysis revealed that 64 provided the optimal balance between gradient stability and noise-driven exploration, while optimizer comparisons showed AdamW clearly outperforming SGD variants by achieving faster, smoother convergence and balanced predictions. A learning rate of 0.005 combined with early stopping yielded the most robust generalization, confirming the value of moderate adaptivity in learning dynamics.

- **Task 6 – Overfitting and Regularization**

The final task assessed regularization strategies. The new baseline deep model already demonstrated strong generalization (accuracy 96.4%). Light weight decay (10^{-4}) provided the best trade-off between accuracy and stability, while stronger methods—Batch Normalization and Dropout—caused underfitting and class collapse. These results indicated that excessive regularization can degrade performance on tabular intrusion data, where the learning signal is already sparse and noise-sensitive.

Final Reflection

The findings demonstrate that integrating bias mitigation, adaptive optimization, and balanced loss design can substantially enhance intrusion detection reliability. The optimized FFNN achieved stable convergence, high interpretability, and strong minority-class recall without excessive complexity. Such improvements contribute directly to the development of more resilient Intrusion Detection Systems (IDS), capable of recognizing both prevalent and rare attack types in real-world network environments.