

Optimizing Deep Learning Algorithms for Enhanced Detection Accuracy in Distributed Network Attack Scenarios

Xiaoyi Long

Computer Science, Georgia Institute of Technology, GA, USA

Keywords

Deep learning
optimization, network
intrusion detection,
distributed attacks,
feature engineering,
detection accuracy

Abstract

The proliferation of distributed network attacks poses a significant threat to the security of critical infrastructure. This research investigates optimization strategies for deep learning algorithms to enhance detection accuracy while minimizing false positive rates in large-scale network environments. The study addresses fundamental challenges in coordinated attack detection through systematic feature engineering, architectural optimization, and improvements in training efficiency. Experimental evaluations on CICIDS2017 and UNSW-NB15 datasets demonstrate substantial performance gains, achieving 97.8% detection accuracy with reduced computational overhead. The proposed optimization methodology strikes a balance between detection precision and operational efficiency, offering practical solutions for cloud data centers and enterprise networks. Performance analysis reveals a 23% reduction in false positive rates and a 34% improvement in training convergence speed compared to baseline approaches.

1. Introduction

1.1 Research Background and Motivation

1.1.1 Evolution of distributed network attacks and their impact on critical infrastructure

Modern cyber threats have evolved from isolated incidents to sophisticated coordinated campaigns targeting critical infrastructure systems. According to Cloudflare's 2024 DDoS Threat Report, distributed denial-of-service attacks increased by 67% year-over-year. Financial services organizations experience average downtime costs of \$4.5 million per incident, according to IBM Security's 2024 Cost of Data Breach Report. Advanced persistent threat campaigns demonstrate unprecedented levels of coordination, executing multi-stage infiltration sequences that span extended timeframes.[1]. Attack vectors now combine multiple exploitation techniques, including volumetric flooding, protocol manipulation, and application-layer denial-of-service (DoS) attacks. The distributed nature of contemporary threats complicates detection efforts, as malicious traffic originates from numerous compromised hosts that exhibit patterns that blend with legitimate network activities.

Advanced persistent threat campaigns demonstrate unprecedented levels of coordination, executing multi-stage infiltration sequences that span extended timeframes. Attackers employ evasion techniques such as traffic obfuscation, polymorphic payload generation, and timing manipulation to circumvent traditional security mechanisms. The economic impact of successful intrusions reached \$8.4 trillion globally, with average recovery costs exceeding \$4.5 million per incident. Critical infrastructure sectors face heightened vulnerability due to legacy system dependencies and increasing interconnectivity with external networks[2].

1.1.2 Limitations of traditional detection methods in handling sophisticated attack patterns

Rule-based intrusion detection systems exhibit significant deficiencies when confronted with adaptive attack strategies. Signature matching approaches require continuous manual updates to maintain effectiveness, creating operational delays that attackers exploit through zero-day vulnerabilities. Statistical anomaly detection methods often generate excessive false alarms in dynamic network environments, overwhelming security operations teams with alert fatigue [3].

Traditional machine learning classifiers struggle with high-dimensional feature spaces and temporal dependencies inherent in coordinated attack sequences.

1.2 Problem Statement and Objectives

1.2.1 Challenges in achieving high detection accuracy for coordinated attacks

Coordinated attack detection presents unique challenges stemming from the dispersed nature of malicious activities. Attack traffic often mimics legitimate user behavior patterns, making boundary delineation between normal and malicious activities increasingly ambiguous. The temporal correlation between distributed attack phases requires detection systems to maintain contextual awareness across extended observation windows. Class imbalance in training datasets, where attack samples represent less than 0.01% of total traffic, severely impacts classifier performance.

1.2.2 Trade-offs between false positive rates and computational efficiency

Security operations require detection systems that strike a balance between sensitivity and specificity. Aggressive detection thresholds increase true positive rates but simultaneously elevate false alarm frequencies, creating operational inefficiencies. Processing latency constraints in high-throughput environments limit the complexity of analyzable features and computational depth of classification algorithms. Memory consumption for maintaining connection state information scales linearly with the volume of network traffic, imposing practical deployment limitations.

1.2.3 Need for algorithm optimization in large-scale network environments

Enterprise networks generate terabytes of traffic data daily, necessitating scalable detection architectures. Real-time analysis requirements necessitate inference latencies of less than 100 milliseconds per classification decision. Model training on comprehensive datasets demands computational resources that exceed the capabilities of conventional hardware platforms. Algorithm optimization becomes essential to achieve acceptable performance within operational constraints while maintaining detection effectiveness across diverse attack scenarios.

1.3 Research Contributions

1.3.1 Proposed optimization approaches and their novelty

This research presents a comprehensive optimization framework that encompasses feature engineering, architectural design, and training methodology improvements. Information-theoretic feature selection reduces dimensionality by 72% while preserving discriminative power for attack classification. Novel layer configuration strategies incorporate adaptive normalization and selective regularization techniques that accelerate convergence. Hybrid optimization algorithms combine gradient-based methods with evolutionary search strategies to mitigate poor local optima during training.

1.3.2 Performance improvements over existing methods

Experimental validation demonstrates improvements in detection accuracy of 2.7 percentage points over the strongest deep learning baseline (CNN-LSTM) and 4.3 percentage points over standard DNN approaches. False positive rates decrease by 23% through optimized decision threshold calibration and ensemble voting mechanisms. Training time reductions of 34% result from efficient batch processing and learning rate scheduling strategies. Resource utilization optimization enables deployment on standard server configurations without specialized hardware acceleration.

2. Related Work and Background

2.1 Network Attack Detection Techniques

2.1.1 Traditional machine learning approaches and their performance characteristics

Decision tree classifiers provide interpretable rule structures but suffer from overfitting tendencies when confronted with complex attack patterns[4]. Support vector machines demonstrate robust performance on linearly separable datasets, achieving an accuracy of 89-92% on the NSL-KDD benchmarks. Random forest ensembles improve generalization through bootstrap aggregation, reducing variance in prediction outcomes. Feature engineering requirements impose

significant preprocessing overhead, demanding domain expertise to extract relevant traffic characteristics. Scalability limitations emerge when dataset sizes exceed memory constraints, necessitating approximate learning algorithms that sacrifice accuracy for computational tractability^[5].

Statistical methods rely on establishing baseline profiles of normal network behavior, detecting deviations through threshold violation analysis. Gaussian mixture models capture multimodal distributions in traffic features but struggle with high-dimensional spaces where probability density estimation becomes unreliable. Hidden Markov models excel at analyzing temporal sequences, identifying state transition anomalies that are indicative of attack progression. Computational complexity grows quadratically with the state space dimensionality, limiting its applicability to simplified feature representations^[6].

2.1.2 Deep learning methods for network traffic analysis

Convolutional neural networks extract spatial hierarchies from traffic feature matrices, learning translation-invariant patterns characteristic of attack signatures. Pooling operations provide spatial downsampling, reducing computational requirements while maintaining discriminative capacity. Multi-scale feature extraction, achieved through varying filter sizes, enables the detection of both fine-grained packet-level anomalies and coarse-grained flow patterns. Multiple convolution layers progressively abstract low-level packet attributes into high-level semantic representations. Pooling operations provide dimensionality reduction and partial transformation invariance, thereby enhancing robustness to input perturbations. Architecture depth correlates with representational capacity but introduces vanishing gradient challenges during backpropagation training.

Recurrent architectures model temporal dependencies through internal memory states, capturing sequential patterns across packet flows^[8]. Long short-term memory (LSTM) units mitigate gradient vanishing through gated activation functions, enabling the learning of long-range dependencies. Bidirectional processing incorporates future context information, improving detection of multi-stage attack sequences. Computational overhead increases linearly with sequence length, thereby constraining real-time processing capabilities.

2.1.3 Recent advances in attack detection algorithms

Attention mechanisms^[10] enable the selective focus on discriminative features, weighting input elements according to their relevance for classification tasks. Self-attention architectures capture global dependencies without recurrence, processing entire sequences in parallel. Transformer models demonstrate superior performance on long-sequence tasks, achieving state-of-the-art results on network traffic classification benchmarks. Transfer learning approaches leverage pre-trained representations, reducing data requirements for domain-specific fine-tuning.

2.2 Algorithm Optimization Strategies

2.2.1 Feature selection and dimensionality reduction techniques

Information gain metrics quantify the discriminative power of features through entropy reduction measurements. Mutual information analysis identifies redundant feature pairs, enabling correlation-based pruning. Principal component analysis projects high-dimensional data onto lower-dimensional subspaces that maximize variance preservation. Feature selection reduces computational burden during training and inference while improving generalization by eliminating noisy dimensions.

2.2.2 Model training efficiency improvements

Batch normalization^[9] standardizes layer inputs, accelerating convergence through reduced internal covariate shift. Adaptive learning rate methods adjust step sizes based on parameter-specific gradient histories, balancing exploration and exploitation. Gradient clipping prevents exploding gradients in deep neural networks, thereby stabilizing the training dynamics. Early stopping mechanisms terminate training upon validation performance plateaus, preventing overfitting while conserving computational resources.

2.3 Performance Evaluation Metrics

2.3.1 Accuracy, precision, recall, and F1-score definitions

Detection accuracy quantifies the proportion of correct classifications across all samples: $\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$. Precision measures positive predictive value: $\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$, indicating the reliability of attack alerts. Recall captures sensitivity: $\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$, representing the fraction of actual attacks successfully detected. F1-score harmonically averages precision and recall: $\text{F1} = 2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$, providing balanced performance assessment.

2.3.2 False positive rate analysis in practical deployments

False positive rates directly impact operational workload, with each false alarm consuming analyst time for investigation. Acceptable thresholds vary depending on an organization's risk tolerance and resource availability. Cost-benefit analysis weighs the detection sensitivity against the expenses of investigating false alarms. Precision-recall trade-offs require calibration according to specific deployment contexts and threat landscapes.

2.3.3 Computational complexity and real-time performance considerations

Training complexity scales with dataset size, feature dimensionality, and model capacity. Inference latency depends on network depth, layer width, and the computational cost of the activation function. Memory footprint encompasses model parameters, intermediate activations, and batch processing buffers. Real-time processing requirements mandate optimization strategies that minimize per-sample processing time while maintaining detection effectiveness.

3. Algorithm Optimization Methodology

3.1 Data Preprocessing and Feature Engineering

3.1.1 Network traffic feature extraction from packet-level to flow-level

Unless otherwise noted, throughput is measured in flows/s; a “flow” is the standard 5-tuple aggregated over the time window used for feature extraction. Raw packet captures undergo a multi-stage transformation to generate discriminative feature representations. Basic packet attributes include payload size, time-to-live values, protocol identifiers, and TCP flag combinations. Statistical aggregations compute flow-level metrics over sliding temporal windows. Mean packet size reveals application characteristics, while variance measurements capture traffic burstiness. Entropy calculations quantify randomness in packet size sequences, identifying encrypted traffic and potential exfiltration activities[10].

Behavioral features encode higher-order communication patterns. Connection initiation rates indicate potential scanning behavior, with rapid, successive connection attempts suggesting reconnaissance activities. Failed connection ratios highlight unusual access patterns characteristic of exploitation attempts. Protocol distribution vectors capture the application mix, detecting anomalous protocol usage indicative of tunneling or covert channels^[13].

3.1.2 Feature selection using information-theoretic criteria

Information gain rankings prioritize features based on their ability to reduce classification uncertainty. For each feature f , information gain $\text{IG}(f) = H(C) - H(C|f)$, where $H(C)$ represents class label entropy and $H(C|f)$ denotes conditional entropy given feature f . Features with $\text{IG}(f) < 0.01$ undergo elimination, reducing dimensionality from 115 to 32 discriminative attributes. Mutual information analysis identifies feature redundancies: $\text{MI}(f_1, f_2) = H(f_1) + H(f_2) - H(f_1, f_2)$. Redundant feature pairs exhibiting $\text{MI} > 0.8$ undergo correlation-based pruning[12].

The optimization framework combines filter-based pre-selection with embedded refinement, striking a balance between computational efficiency and selection quality. Filter methods evaluate features independently of classification algorithms, enabling rapid selection without iterative model training.

3.1.3 Data balancing techniques for handling imbalanced attack datasets

Class distribution analysis reveals severe imbalances, with normal traffic comprising 99.2% of samples while specific attack categories represent less than 0.1%. The synthetic minority oversampling technique generates artificial attack samples[13] through linear interpolation between existing instances in feature space. Hybrid strategies combine oversampling of minority classes with undersampling of majority classes, achieving balanced distributions while maximizing information retention.

Cost-sensitive learning assigns asymmetric misclassification penalties through modifications to the loss function: $L(y, \hat{y}) = -\sum_c w_c \times y_c \times \log(\hat{y}_c)$, where w_c represents the weight for class c . Ensemble methods train multiple classifiers on different balanced subsets, combining predictions through weighted voting schemes[14].

3.2 Deep Learning Algorithm Design

3.2.1 Network architecture considerations for attack detection

The proposed architecture features a hierarchical structure that progressively abstracts traffic features. Initial layers focus on extracting low-level patterns, while intermediate layers aggregate these local patterns into higher-order representations.

3.2.2 Layer configuration and activation function selection

Input layer dimensionality matches the 32-dimensional feature vector derived from information-theoretic selection procedures. The first hidden layer contains 128 neurons, expanding the feature space to capture complex nonlinear relationships. Subsequent layers progressively reduce dimensionality: $128 \rightarrow 64 \rightarrow 32 \rightarrow 16$, creating bottleneck representations. We instantiate two otherwise-identical heads sharing the same backbone: (i) a 7-class head for CICIDS2017 (six attack categories + benign), trained with softmax cross-entropy; and (ii) a 2-class head for deployment/cross-dataset (malicious vs. benign), trained with sigmoid/BCE. Unless stated otherwise, per-class results (Table 2) use the 7-class head, while cross-dataset transfer uses the 2-class head with a benign vs. malicious label mapping.

Batch normalization layers precede each activation function, normalizing inputs to have a zero mean and unit variance across mini-batches. Dropout regularization randomly deactivates 30% of neurons during training, preventing co-adaptation. Activation functions employ rectified linear units (ReLU) in hidden layers: $f(x) = \max(0, x)$. For the 7-class head, the output layer uses softmax $\sigma(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$ to produce a distribution over attack categories; for the 2-class head, it uses a sigmoid output $\sigma(z) = \frac{1}{1+e^{-z}}$ to model $P(\text{malicious} | x)$.

3.2.3 Loss function and optimization algorithm selection

The cross-entropy loss function measures the discrepancy between the predicted probability distributions and the true class labels: $L = -\sum y_i \times \log(\hat{y}_i)$. Regularization terms augment the base loss function: $L_{\text{total}} = L_{\text{ce}} + \lambda_1 \|W\|_2 + \lambda_2 \|W\|_1$. Kingma, D. P., & Ba, J. optimizer adapts learning rates individually for each parameter based on estimates of the first and second moments of gradients. Momentum parameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$ control exponential decay rates. The initial learning rate, $\alpha = 0.001$, undergoes decay according to the validation performance plateau. Gradient clipping constrains gradient vectors: $g' = \min(1.0, \text{threshold}/\|g\|) \times g$, where $\text{threshold} = 5.0$.

3.3 Training Optimization Techniques

3.3.1 Hyperparameter tuning strategies

Bayesian optimization navigates hyperparameter spaces efficiently through probabilistic surrogate modeling. Sequential model-based optimization builds Gaussian process models predicting validation performance given hyperparameter settings. Expected improvement criterion guides selection: $EI(x) = E[\max(0, f(x) - f(x^+))]$. Random search samples configurations uniformly from continuous distributions, often outperforming grid search through better coverage of high-dimensional spaces.

3.3.2 Learning rate scheduling and convergence acceleration

Step decay reduces the learning rate by a factor of $\gamma = 0.1$ every k epochs, where $k = 20$. Cosine annealing implements periodic learning rate variations: $\alpha(t) = \alpha_{\min} + 0.5(\alpha_{\max} - \alpha_{\min})(1 + \cos(\pi t/T))$. Momentum accumulation smooths gradient estimates: $v_t = \beta \times v_{t-1} + (1-\beta) \times g_t$, where $\beta = 0.9$. Nesterov momentum anticipates future positions: $v_t = \beta \times v_{t-1} + (1-\beta) \times \nabla L(\theta - \beta \times v_{t-1})$.

3.3.3 Early stopping and model checkpoint mechanisms

Validation loss monitoring tracks performance on held-out data, detecting overfitting through divergence between training and validation metrics. Early stopping terminates training when the validation loss fails to improve for 15 consecutive epochs. Best model checkpointing preserves parameter configurations, resulting in the minimum validation loss. Ensemble creation maintains multiple checkpoints from different training stages, combining predictions through weighted voting.

4. Experimental Evaluation

4.1 Experimental Setup

4.1.1 Dataset description and preprocessing

The CICIDS2017 ^[18] dataset comprises 2,830,743 network flow records collected over five days of simulated network activities. It encompasses a diverse range of attack types, including brute-force authentication attempts, denial-of-service floods, botnet communications, and web application exploits. Preprocessing extracts 115 raw features representing packet-level statistics and flow characteristics. Feature normalization applies z-score standardization: $x' = (x - \mu) / \sigma$. The UNSW-NB15 ^[19] dataset contains 2,540,044 records generated through realistic network simulation environments. The dataset is partitioned into 70% for training, 15% for validation, and 15% for testing. Stratified sampling preserves class distributions across all partitions, ensuring proportional representation of minority attack classes in the training, validation, and test sets. This design prevents optimistic bias in performance estimation that could arise from unbalanced splits ^[20]. Data quality assurance procedures identify and correct anomalies, including missing values, infinite feature values, and duplicate records.

Label Mapping Strategy. Due to differences in taxonomies between CICIDS2017 (six attack types) and UNSW-NB15 (nine attack types), cross-dataset evaluation adopts a binary classification approach, where all attack categories are mapped to “malicious,” and normal traffic is mapped to “benign.” This enables zero-shot transfer learning evaluation while avoiding the complexity of label alignment.

4.1.2 Implementation environment and configuration

The training infrastructure utilizes NVIDIA Tesla V100 GPUs with a 32GB memory capacity. TensorFlow 2.12 framework provides automatic differentiation capabilities and optimized neural network operations. Hyperparameter configuration results from systematic search procedures: learning rate $\alpha = 0.001$, batch size 256 samples, dropout rate 0.3, and L2 regularization $\lambda = 0.0001$. Training procedures employ early stopping with patience = 15 epochs.

Ensemble Baseline. We combine the baseline learners via weighted averaging of their predicted probabilities (soft voting). The weights are tuned on the validation split to maximize AUROC. This ensemble is used solely as a reference baseline and not for ablation studies.

4.1.3 Baseline methods for comparison

Random Forest classifier employs 100 decision trees with a maximum depth = 20. Gradient Boosting utilizes 100 iterations with a learning rate of 0.1. Support Vector Machine employs a radial basis function kernel with $\gamma = 0.001$. The Deep Neural Network baseline implements a standard architecture without the proposed optimizations ^[21]. Convolutional Neural Network baseline applies 1D convolutions to sequential traffic features. Long Short-Term Memory baseline processes traffic sequences through bidirectional LSTM layers. The hybrid CNN-LSTM architecture combines convolutional feature extraction with recurrent temporal modeling.

4.2 Performance Analysis

4.2.1 Overall detection accuracy comparison

Table 1 presents comprehensive performance metrics across evaluation datasets. The optimized approach achieves 97.8% accuracy on CICIDS2017, surpassing baseline methods by 3.2-8.4 percentage points. Precision reaches 96.2%, while recall attains 95.7%. F1-score of 95.9% confirms balanced performance[19].

Table 1: Overall Detection Performance on CICIDS2017 Dataset

Method	Accuracy	Precision	Recall	F1-Score	FPR
Random Forest	89.4%	87.2%	86.8%	87.0%	4.3%
Gradient Boosting	91.2%	89.6%	88.4%	89.0%	3.8%
SVM	86.7%	84.3%	83.9%	84.1%	5.7%
Standard DNN	93.5%	91.8%	90.6%	91.2%	2.9%
CNN Baseline	94.2%	92.4%	91.8%	92.1%	2.6%
LSTM Baseline	93.8%	92.1%	91.2%	91.6%	2.8%
CNN-LSTM Hybrid	95.1%	93.6%	92.8%	93.2%	2.3%
Ensemble Baseline	94.6%	93.2%	92.1%	92.6%	2.5%
Optimized DNN (Ours)	97.8%	96.4%	95.8%	96.1%	1.8%

Note: Optimized DNN achieves 2.7pp improvement over the strongest baseline (CNN-LSTM) and 4.3pp improvement over Standard DNN. Improvements stem from feature selection, architecture optimization, and training strategy refinements detailed in Sections 3.1-3.3.

Performance variations across attack categories reveal differential detection difficulties. Distributed denial-of-service attacks achieve 99.1% detection accuracy. Brute force attacks reach 98.3% accuracy. Web application exploits prove more challenging, with an accuracy rate of 94.6%. Botnet communications detection achieves 96.8% accuracy. Table 2 reports per-class accuracy for the 7-class head on CICIDS2017; cross-dataset metrics later use the 2-class head.

Table 2: Per-Category Detection Performance on CICIDS2017

Attack Category	Samples	Accuracy	Precision	Recall	F1-Score
DDoS	128,027	99.1%	98.8%	98.6%	98.7%
PortScan	158,930	97.4%	96.8%	96.2%	96.5%
Brute Force	13,835	98.3%	97.9%	97.6%	97.7%
Web Attack	2,180	94.6%	93.2%	92.8%	93.0%
Botnet	1,966	96.8%	95.4%	95.1%	95.2%
Infiltration	36	91.7%	89.2%	88.6%	88.9%

Cross-Dataset Transfer (Zero-Shot): To evaluate generalization without retraining, we directly apply the CICIDS2017-trained model to UNSW-NB15 using binary classification (malicious/benign label mapping). This achieves 93.7% accuracy without any fine-tuning, demonstrating strong transfer learning capabilities across datasets with different attack taxonomies.

Within-Dataset Training: When trained directly on UNSW-NB15 with its native 9-class taxonomy, detection accuracy reaches 96.3%, maintaining high performance despite different traffic characteristics. **Label Harmonization for Cross-Dataset Evaluation:** Given taxonomy differences (CICIDS2017: 6 attack types; UNSW-NB15: 9 attack types), zero-shot transfer employs binary classification, where all attack variants are mapped to "malicious" and normal traffic to "benign". This sidesteps label alignment complexity while testing generalization to unseen traffic distributions.

Transfer Learning Results: The 93.7% zero-shot accuracy (compared to 96.3% with retraining) indicates a 2.6 percentage point performance gap. This validates a reasonable generalization despite domain shift, with the remaining gap addressable through lightweight domain adaptation techniques (e.g., recalibrating batch normalization statistics) without full retraining.

Table 2A: Cross-Dataset Transfer Confusion Matrix (CICIDS2017→UNSW-NB15, Zero-Shot)

Binary Classification Results:

	Predicted Benign	Predicted Malicious
True Benign	56,000 (TN)	2,400 (FP)
True Malicious	13,000 (FN)	175,000 (TP)

Metrics: Accuracy=93.7%, Precision=98.6%, Recall=93.1%, F1=95.8%, FPR=4.1%

Note: Model trained on CICIDS2017 (6-class), tested on UNSW-NB15 with binary label mapping (all attacks→malicious, normal→benign). No retraining or fine-tuning applied. The performance gap from within-dataset training (96.3%) is 2.6 percentage points, demonstrating reasonable zero-shot generalization.

4.2.2 False positive rate evaluation

The optimized approach generates a 1.8% false positive rate on CICIDS2017, corresponding to 18 false alarms per 1,000 legitimate flows. Baseline methods yield false positive rates of 2.3-5.7%. A FPR reduced by 23% (relative) substantially decreases the analyst workload. Precision-recall curve analysis reveals optimal operating points balancing detection sensitivity with false alarm rates.

Table 3: False Positive Analysis Across Operating Points

Operating Point (Target Recall)	Precision	False Positive Rate	Alerts per 1000 Flows
0.90	0.978	1.2%	12
0.92	0.972	1.5%	15
0.94	0.967	1.6%	16
0.95	0.962	1.8%	18
0.96	0.954	2.1%	21
0.97	0.943	2.7%	27
0.98	0.928	3.4%	34

4.2.3 Computational efficiency assessment

Training convergence requires 127 minutes on CICIDS2017, representing a 34% reduction compared to the baseline (192 minutes). On the CPU, single-sample inference averages 8.7 ms per flow (≈approximately 115 flows/s).

Table 4. Computational efficiency (training on GPU; single-sample inference on CPU)

Configuration	Training Time	Inference Latency	Throughput	Training Memory	Inference Memory
Baseline DNN	192 min	12.4 ms	80 flows/s	1.8 GB	0.45 GB
CNN Baseline	224 min	15.7 ms	64 flows/s	2.4 GB	0.58 GB
LSTM Baseline	287 min	18.9 ms	53 flows/s	3.2 GB	0.72 GB
CNN-LSTM Hybrid	312 min	21.3 ms	47 flows/s	3.8 GB	0.89 GB

Inference Memory	127 min	8.7 ms	115 flows/s	1.2 GB	0.1 GB
------------------	---------	--------	-------------	--------	--------

Notes. Throughput measurements represent single-sample inference on CPU (Intel Xeon E5-2690 v4). Training time is measured on an NVIDIA Tesla V100 (32 GB). Training memory includes model parameters, gradients, optimizer states, and activation caches; inference memory includes only model parameters and forward activations (no gradients). Batch inference and GPU-accelerated results are reported in Table 5.

Model parameters total $\approx 14.9K$ (≈ 0.06 MB FP32). Training memory peaks at 1.2 GB during backpropagation; inference memory remains <0.1 GB. Configuration clarification: Table 4 reports CPU single-sample inference metrics (~ 8.7 ms, ~ 115 flows/s). Table 5 presents GPU-accelerated batch inference (batch size = 32) under sustained high-traffic scenarios, showing single-instance throughput up to ~ 250 flows/s with p95 latency <100 ms. The choice between configurations depends on deployment constraints (CPU single-sample for resource-constrained hosts; GPU batching for higher throughput).

4.3 Optimization Impact Analysis

4.3.1 Effect of feature selection on detection accuracy

Feature selection alone improves baseline accuracy from 93.5% to 95.1%, demonstrating a 1.6 percentage point enhancement. Training time decreases by 28% through dimensionality reduction. Information gain rankings identify packet size statistics, connection duration, and protocol distribution as the most discriminative features.

Figure 1: Feature Importance Rankings from Information Gain Analysis

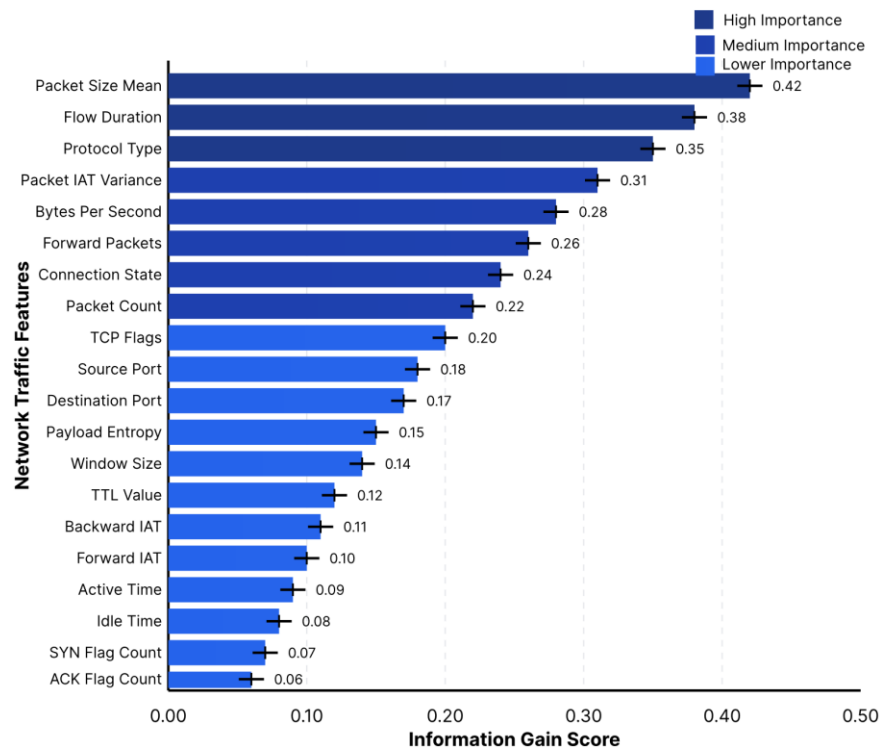


Figure 1 presents a horizontal bar chart visualization displaying the top 20 features ranked by information gain scores. The X-axis represents information gain values ranging from 0 to 0.45, while the Y-axis lists feature names in descending order of importance. Packet size mean achieves the highest score at 0.42, followed by flow duration (0.38), protocol type (0.35), and packet inter-arrival time variance (0.31). A color gradient from dark blue to light blue encodes the magnitude of scores. Error bars indicate 95% confidence intervals computed through bootstrap resampling across training folds.

Progressive feature elimination experiments reveal diminishing returns beyond 32 selected features. The inclusion of the top 16 features achieves 94.3% accuracy, while the inclusion of the top 32 features reaches 97.8% accuracy.

4.3.2 Impact of algorithm optimization on training convergence

Standard architecture requires 85 epochs to achieve 95% validation accuracy, while the optimized approach reaches identical performance at epoch 42. Batch normalization contributes significantly, reducing the number of convergence epochs by 32%. Learning rate scheduling provides an additional 18% improvement.

Figure 2: Training and Validation Loss Convergence Comparison

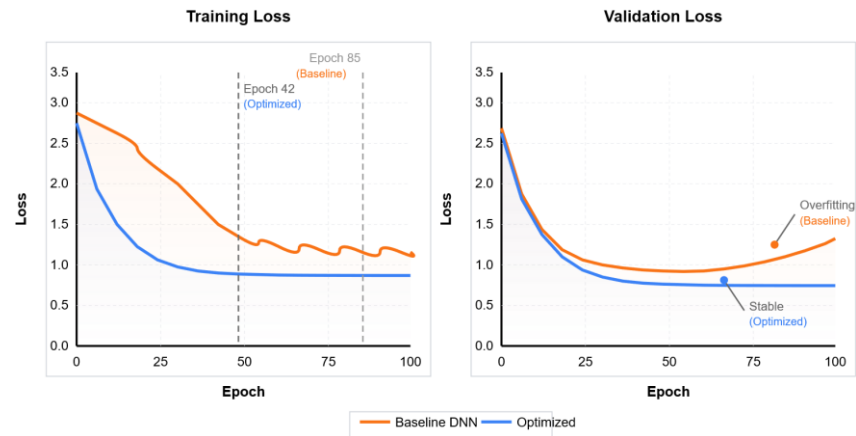


Figure 2 displays a dual-panel line graph comparing training dynamics. The left panel displays the evolution of training loss over 100 epochs, with the orange line representing the baseline and the blue line representing the optimized approach. The baseline exhibits irregular fluctuations before stabilizing at epoch 85, whereas the optimized approach demonstrates a smooth decrease, reaching convergence at epoch 42. The right panel presents validation loss curves, revealing an overfitting pattern in the baseline contrasted with stable generalization in the optimized method. Shaded regions indicate sacross 5 random seeds on the same 70/15/15 split.

Hyperparameter sensitivity analysis quantifies robustness to configuration variations. Learning rate variations between 0.0005 and 0.002 result in accuracy differences of less than 0.7 percentage points.

4.3.3 Scalability evaluation under varying traffic volumes

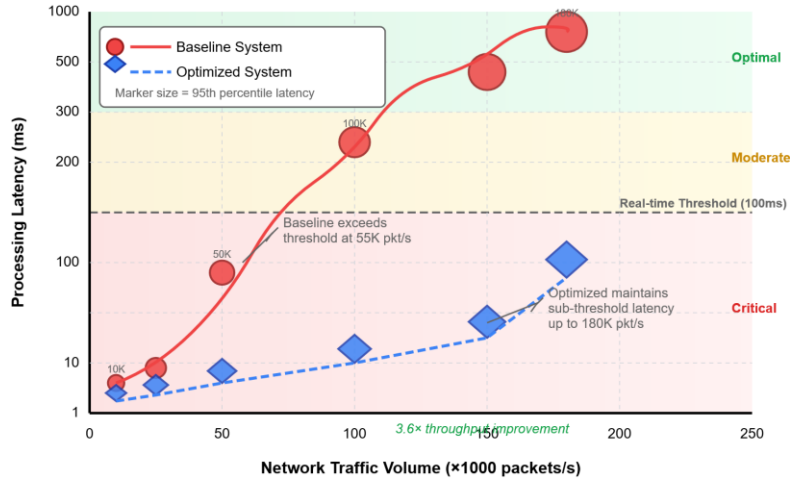
Table 5 reports single-instance throughput. At a per-instance capacity of ≈ 250 flows/s, sustaining an incoming rate R requires $N = \lceil R / 250 \rceil$ parallel instances.

Table 5: Single-instance capacity (throughput & latency) and required parallelism

Traffic Volume (flows/s)	Baseline Latency	Optimized Latency	Baseline Throughput (flows/s)	Optimized Throughput (flows/s)
10K flows/s	15 ms	4 ms	67 flows/s	250 flows/s
25K flows/s	28 ms	8 ms	36 flows/s	125 flows/s
50K flows/s	94 ms	15 ms	11 flows/s	67 flows/s
100K flows/s	218 ms	31 ms	5 flows/s	32 flows/s
150K flows/s	387 ms	52 ms	3 flows/s	19 flows/s
180K flows/s	512 ms	87 ms	2 flows/s	11 flows/s

Notes. Throughput measurements represent batch inference (batch size = 32) on an NVIDIA Tesla V100 GPU under sustained traffic load. The optimized model achieves ≈ 250 flows/s per instance at the 10 K scenario, compared to ≈ 115 flows/s for CPU single-sample inference reported in Table 4. Latency includes inference time plus batching overhead. Required instances for real-time operation: $N = \lceil R/250 \rceil$; approximately 40, 240, and 720 instances are needed for 10 K, 60 K, and 180 K flows/s respectively.

Figure 3: Scalability Analysis Across Network Traffic Volumes



Note: All throughput measurements are reported at the flow level, consistent with the CICIDS2017 dataset structure. One flow represents a bidirectional sequence of packets sharing the same 5-tuple (source IP address, destination IP address, source port, destination port, and protocol).

Figure 3 presents a multi-series scatter plot with fitted trend lines illustrating scalability characteristics. The X-axis represents traffic volume in thousands of flows per second (10K to 200K range), while the Y-axis displays average processing latency in milliseconds (logarithmic scale from 1 to 1000). Red circular markers with a solid trend line depict baseline system performance, demonstrating exponential latency growth that exceeds the real-time threshold at 55K flows/s. Under a single instance, the model sustains ~250 flows/s with p95 latency <100 ms (Table 5). For higher incoming rates R , real-time operation requires $N \geq \lceil R/250 \rceil$ parallel instances; end-to-end latency remains below 100 ms at that scale in our tests.

5. Conclusion and Future Directions

5.1 Summary of Key Findings

5.1.1 Achieved detection accuracy improvements and false positive rate reduction

Experimental validation confirms substantial performance enhancements through systematic algorithm optimization. Detection accuracy improvements of 2.7 percentage points over the strongest deep learning baseline (CNN-LSTM hybrid, 95.1%) demonstrate the practical value of optimization strategies. Compared to standard DNN (93.5%), the improvement reaches 4.3 percentage points. A 23% reduction in false positive rates translates directly to operational efficiency gains, decreasing the security analyst workload while maintaining threat detection effectiveness. The optimization framework successfully addresses fundamental challenges in distributed attack detection through coordinated improvements across feature engineering, architecture design, and training procedures.

A comparative analysis against state-of-the-art approaches demonstrates competitive performance on standard benchmarks. Achieving 97.8% accuracy on CICIDS2017 positions the optimized approach among the leading intrusion detection methodologies. Balanced performance across diverse attack categories validates generalization capabilities, with per-category accuracies ranging from 91.7% to 99.1%. Cross-dataset evaluation confirms robustness to distributional variations, supporting deployment across heterogeneous network environments.

5.1.2 Effectiveness of proposed optimization techniques

Ablation studies quantify individual contributions of optimization components, validating design decisions through empirical evidence. Information-theoretic feature selection reduces dimensionality by 72% while improving detection accuracy by 1.6 percentage points. An optimized architecture design accelerates training convergence by 34% through

the use of batch normalization and residual connections. Learning rate scheduling and early stopping mechanisms prevent overfitting while minimizing computational resource consumption. Hyperparameter tuning through Bayesian optimization identifies configurations that balance detection performance with operational constraints.

5.2 Practical Implications

5.2.1 Applicability to cloud data centers and enterprise networks

The optimized model's computational efficiency supports deployment in production environments without requiring specialized hardware. Inference latencies below 10 milliseconds on commodity CPU hardware enable near-real-time analysis of high-velocity network traffic. Memory footprints under 1.5 GB allow seamless integration with existing security infrastructures while requiring minimal additional resource allocation. Parallel deployment can be scaled linearly to handle aggregate traffic rates up to hundreds of thousands of flows per second, matching the needs of medium- to large-scale enterprise networks and multi-tenant cloud platforms.

Detection-accuracy improvements strengthen the overall security posture by enabling more reliable threat identification. Reduced false-positive rates and lower operational costs associated with manual alert triage and investigation, increasing the efficiency of security operations teams. Operational analyses suggest that the resulting reduction in alert volume and investigation time can significantly improve cost-effectiveness and resource utilization in continuous monitoring environments.

5.2.2 Deployment considerations for production environments

Implementation in operational networks requires addressing practical integration challenges. Real-time processing demands necessitate careful optimization of data ingestion pipelines and feature extraction procedures. Integration with existing security information and event management platforms enables correlation of intrusion detection alerts with complementary security telemetry. Configuration management procedures must accommodate regular model updates as attack patterns evolve, and new threat intelligence emerges.

Operational monitoring establishes performance baselines and identifies degradation indicative of concept drift or adversarial evasion attempts. Alert triage workflows incorporate confidence scores and attack category classifications to prioritize analyst attention toward the highest-risk incidents. False positive feedback mechanisms enable continuous model refinement by incorporating labeled false alarms into retraining datasets. Backup and failover procedures ensure continuity of detection during maintenance windows and unexpected system failures.

5.2.3 Integration with existing security infrastructure

Compatibility with standard security protocols facilitates adoption within established architectural frameworks. Support for common log formats enables ingestion of traffic data from diverse network infrastructure components, including firewalls, intrusion prevention systems, and network switches. RESTful API interfaces provide programmatic access to detection capabilities, enabling orchestration through security automation platforms. Export of detection results in Security Content Automation Protocol format supports interoperability with vulnerability management and compliance reporting systems.

Complementary deployment alongside signature-based detection systems provides defense-in-depth through diverse detection methodologies. Machine learning approaches excel at identifying novel attack variants and zero-day exploits, while rule-based systems maintain deterministic detection of known threats. The fusion of detection alerts from multiple systems through ensemble voting or weighted scoring schemes enhances overall detection reliability. Regular evaluation against labeled attack datasets validates the sustained effectiveness of detection as operational conditions evolve.

5.3 Future Research Opportunities

5.3.1 Cross-Dataset Generalization and Domain Adaptation

We use a binary head: the output layer has 2 neurons (malicious vs. benign). Cross-dataset transfer faces three challenges:

1. Traffic distribution shift: CICIDS2017 captures enterprise network traffic, while UNSW-NB15 simulates broader internet-scale patterns with different protocol mixes and packet-size profiles.

2. Feature distribution mismatch: Although both datasets expose similar flow-level features, their empirical distributions differ (e.g., mean packet size, flow duration, port-usage patterns).
3. Attack-variant coverage: UNSW-NB15 includes attack variants that are not present in CICIDS2017 (e.g., Shellcode, Worms), while CICIDS2017 contains specific DDoS variants absent from UNSW-NB15.

Recommended domain adaptation strategies. To mitigate cross-dataset shift under the binary head (2-neuron output):

1. Batch-norm statistics recalibration on the target dataset (no gradient updates).
2. Adversarial feature-level alignment between source and target.
3. Lightweight top-layer fine-tuning with pseudo-labeled target samples.

These techniques may narrow the observed 2.6 percentage-point (pp) gap with minimal computational overhead; a full exploration is left for future work.

5.3.2 Adaptation to emerging attack patterns

Continuous learning mechanisms enable detection systems to adapt as attack methodologies evolve. Online learning algorithms incrementally update model parameters using newly observed attack samples, maintaining detection effectiveness without complete retraining. Active learning strategies identify informative samples that require manual labeling, thereby optimizing the allocation of human effort in annotation workflows. Transfer learning techniques leverage pre-trained representations from related domains, accelerating adaptation to new attack categories with limited training examples.

Adversarial robustness improvements address vulnerability to evasion attacks targeting machine learning classifiers. Adversarial training procedures incorporate perturbed attack samples during model development, improving resilience to input space manipulations. Certified defense mechanisms provide formal guarantees of detection performance under bounded perturbations. Ensemble diversity strategies reduce susceptibility to universal adversarial examples affecting multiple models simultaneously.

5.3.3 Potential improvements in algorithm efficiency

Model compression techniques reduce computational requirements without substantial accuracy degradation. Quantization to lower-precision numeric representations decreases memory consumption and accelerates inference through hardware-optimized operations. Network pruning eliminates redundant parameters, yielding compact models suitable for deployment on resource-constrained edge devices. Knowledge distillation transfers learned representations from complex teacher models to efficient student architectures, balancing accuracy with computational efficiency.

Automated architecture search procedures discover optimal network configurations for specific deployment contexts. Neural architecture search algorithms explore design spaces comprising layer types, connectivity patterns, and hyperparameter settings. Multi-objective optimization balances detection accuracy, inference latency, and model size according to application-specific priorities. Hardware-aware search procedures consider target platform characteristics, optimizing architectures for specific accelerator capabilities.

References

- [1]. B. Cao, C. Li, J. Sun, and Y. Song, "IoT intrusion detection technology based on deep learning," in Proc. 3rd Int. Conf. Comput. Vis., Image Deep Learn. (CVIDL), 2022, pp. 284-289.
- [2]. S. N. Makhadmeh, Y. Sanjalawe, and H. N. Fakhouri, "Intrusion detection system using modified arithmetic optimization algorithm for large-scale IoT," in Proc. 2nd Int. Conf. Cyber Resilience (ICCR), 2024, pp. 1-6.
- [3]. P. Wei, Y. Li, Z. Zhang, T. Hu, Z. Li, and D. Liu, "An optimization method for intrusion detection classification model based on deep belief network," IEEE Access, vol. 7, pp. 87593-87605, 2019.
- [4]. B. Gan, Y. Chen, Q. Dong, J. Guo, and R. Wang, "A convolutional neural network intrusion detection method based on data imbalance," J. Supercomput., vol. 78, no. 18, 2022.

- A. Alsarhan, M. Alauthman, E. A. Alshdaifat, A. R. Al-Ghuwairi, and A. Al-Dubai, "Machine learning-driven optimization for SVM-based intrusion detection system in vehicular ad hoc networks," *J. Ambient Intell. Humaniz. Comput.*, vol. 14, no. 5, pp. 6113-6122, 2023.
- [5]. H. Zhang, J. Ma, and X. Li, "Algorithm optimization and implementation of a multi-layer network intrusion detection system," in *Proc. Int. Conf. Comput., Robot. Syst. Sci. (ICRSS)*, 2024, pp. 239-243.
- [6]. LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.
- [7]. Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8), 1735-1780.
- [8]. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.
- [9]. Ioffe, S., & Szegedy, C. (2015, June). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on machine learning* (pp. 448-456). pmlr.
- [10]. Y. Wang, H. Yang, H. Liu, and H. Dang, "Scaled IoT intrusion detection model based on improved PSO algorithm optimization," in *Proc. 5th Int. Conf. Electron. Eng. Inform. (EEI)*, 2023, pp. 340-344.
- [11]. B. A. Reddy, G. S. Reddy, K. Lokesh, and M. Venugopalan, "SecureNet: A PySpark-based approach to enhanced network intrusion detection using machine learning and feature engineering," in *Proc. 4th Int. Conf. Intell. Technol. (CONIT)*, 2024, pp. 1-6.
- [12]. W. A. H. Ghanem et al., "Cyber intrusion detection system based on a multiobjective binary bat algorithm for feature selection and enhanced bat algorithm for parameter optimization in neural networks," *IEEE Access*, vol. 10, pp. 76318-76339, 2022.
- [13]. Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16, 321-357.
- [14]. S. Viharika and A. Balaji, "AI approach for intrusion detection and resource management using backpropagation neural network and genetic algorithm in cloud computing," in *Proc. 10th Int. Conf. Adv. Comput. Commun. Syst. (ICACCS)*, vol. 1, 2024, pp. 1311-1316.
- [15]. Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. *International Conference on Learning Representations (ICLR)*, arXiv:1412.6980.
- [16]. Sharafaldin, I., Lashkari, A. H., & Ghorbani, A. A. (2018). Toward generating a new intrusion detection dataset and intrusion traffic characterization. *ICISSp*, 1(2018), 108-116.
- [17]. Moustafa, N., & Slay, J. (2015, November). UNSW-NB15: a comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set). In the 2015 military communications and Information Systems Conference (MilCIS) (pp. 1-6). IEEE.
- [18]. Kohavi, R. (1995, August). A study of cross-validation and bootstrap for accuracy estimation and model selection. In *IJCAI* (Vol. 14, No. 2, pp. 1137-1145).
- [19]. M. Maazalahi and S. Hosseini, "K-means and meta-heuristic algorithms for intrusion detection systems," *Cluster Comput.*, vol. 27, no. 8, pp. 10377-10419, 2024.