

Bridging the Gap: Interactive Interpretability for Machine Learning-Based Intrusion Detection

Elijah Segura, Joseph Hammock

Department of Computer Science

Central Michigan University

Mount Pleasant, MI, USA

segur1em@cmich.edu, hammo3jm@cmich.edu

Abstract—Machine learning-based Intrusion Detection Systems (IDS) provide strong predictive performance but often at the cost of interpretability, which reduces analyst understanding, trust, and operational usability. This paper presents BridgeIDS, a high-accuracy intrusion detection system that combines an XGBoost classifier with an interactive interpretability layer designed to make model decisions transparent to users in real time. The system is trained on the CSE-CIC-IDS2018 dataset using a small set of engineered network-flow features, and utilizes a hybrid sampling strategy to address severe class imbalance. BridgeIDS integrates SHAP-based explanations, statistical anomaly indicators, and a what-if interface that lets analysts adjust feature values and observe prediction changes with sub-50ms latency. It also provides basic counterfactual guidance, to help users understand how flows could shift from malicious to benign classifications.

Index Terms—Intrusion Detection, Explainable AI, XGBoost, Interactive Visualization, Network Security, SHAP

I. INTRODUCTION

As cyber threats continue to evolve in both scale and sophistication, network security has become an increasingly critical concern. Intrusion Detection Systems (IDS) are essential for identifying malicious activity and preventing attacks, yet traditional signature-based IDS are limited in detecting novel or subtle threats. Machine learning (ML)-based IDS, including ensemble models such as XGBoost [1], have exhibited exceptional detection capabilities by learning complex patterns from historical network traffic. However, these models often act as “black boxes,” providing limited insight into why a particular flow is classified as malicious or benign [2]. This lack of interpretability can reduce trust and understanding in automated alerts, and hinder practical deployment in operational environments.

To address this, we developed BridgeIDS, an ML-based intrusion detection system that combines high-performance classification with interactive interpretability. The system uses an XGBoost classifier trained on the CSE-CIC-IDS2018 dataset [3] to distinguish between six classes: Benign, DoS, DDoS, Brute Force, Web Attack, and Bot/Infiltration. BridgeIDS integrates an interactive web-based dashboard that allows analysts to explore model predictions in real time, understand feature contributions, and conduct “what-if” scenario analyses. Additionally, the system provides counterfactual explanations, offering educational suggestions on feature changes required to reclassify a flow as benign.

The system architecture follows a modular, microservices-like design for scalability and maintainability. The data pipeline handles ingestion, cleaning, preprocessing, and feature engineering, including 12 core flow features (e.g., destination port, flow duration, TCP flags) and 8 derived features (e.g., packet rates, flag ratios, port categories). Class imbalance is addressed through a two-phase strategy: aggressive benign downsampling followed by SMOTE-based oversampling of minority attack classes [4]. Features are standardized to prevent data leakage, and labels are consolidated into six semantic classes for effective training. The XGBoost model is configured with regularization, weighted loss functions, and histogram-based training to optimize both accuracy and efficiency.

A key contribution of BridgeIDS is the interactive interpretability module, which combines statistical context (Z-score analysis), SHAP-based local feature attribution [5], and real-time sensitivity analysis. Analysts can manipulate feature values via sliders and observe prediction shifts in real time (<50ms), revealing relationships learned by the model. Counterfactual “Benign Prescription” helps to further improve insight by showing adjustments needed to cross decision boundaries (e.g., Reduce TCP Window Size to < 24194 (Changes prediction to Benign)).

We evaluated BridgeIDS on a stratified 20% sample of the CSE-CIC-IDS2018 dataset (~12.6 million flows). The system achieved an overall accuracy of 99.96% and weighted F1-score of 0.9996, with 100% recall across all attack types, demonstrating its effectiveness for intrusion detection based on our evaluation sample. Case studies can help illustrate the utility of interactive interpretability: for DoS attacks (Slowloris), excessive SYN flags, high connection duration, and high average time between packets drove the prediction. Reducing the TCP window size, increasing the destination port, or increasing the time of day changed the prediction to benign. For bot traffic, high total packets, bytes, and destination port were key findings. By reducing the time of day (< 11.7) or halving connection duration, the prediction was changed to benign. The system handles extreme data imbalance effectively, maintains low latency (<50ms per prediction), and runs on consumer-grade hardware, positioning it for potential real-time deployment in an educational environment.

BridgeIDS helps to close the gap between high-accuracy

ML classification and actionable interpretability. By providing interactive exploration, semantic insights, and counterfactual reasoning, the system empowers security analysts to detect and understand network intrusions efficiently, while maintaining trust in automated decisions.

II. RELATED WORK

From early anomaly detection methods to sophisticated ensemble approaches, the application of machine learning to intrusion detection has evolved significantly. The KDD Cup 99 dataset was a standard benchmark for many years, but is now regarded as outdated, lacking modern traffic patterns. The CSE-CIC-IDS2018 dataset [3] is a strong option, introduced to address these limitations by providing a realistic cloud-centric dataset with diverse attack scenarios. Recent work has demonstrated the effectiveness of deep learning and ensemble methods on this dataset. Khan et al. [6] utilized deep learning to improve attack detection, addressing class imbalance through data augmentation. Göcs & Johanyák [7] focused on feature selection engineering to identify the most relevant features for efficient classification, a strategy we followed in our selection of 20 high-impact features.

As the demand for machine learning models increases in complexity, the need for interpretability grows. Molnar [2] notes the trade-off between model accuracy and interpretability. He states that high-performing models like XGBoost [1] are often opaque. A systematic review by Mohale and Obagbuwa [8] emphasizes that while XAI is increasingly integrated into IDS, most approaches rely on static, post-hoc explanations. Techniques like SHAP (Shapley Additive Explanations) [5] and LIME have become the standard for local interpretability. However, recent surveys [9], [10] indicate that static explanations alone may not be sufficient for security environments, where analysts need to understand how sensitive model decisions are to changing network conditions, not just why a decision was made.

Current research identifies a gap in actionable, real-time interpretability. Silva et al. [11] and others have explored XAI for IoT and general IDS, but the focus often remains on generating static reports. BridgeIDS advances this by offering an interactive "what-if" analysis capability. This aligns with the challenges identified in recent literature [12], which calls for XAI systems that can support human-in-the-loop (HITL) decision-making. BridgeIDS allows analysts to manipulate feature values and observe the immediate impact on the model's prediction. It moves beyond static SHAP feature importance reports to provide a functional, interactive understanding of the model's decision-making.

Previous ML-based IDS research has explored various algorithms. Random Forest achieves ~95% accuracy on CSE-CIC-IDS2018 but suffers from slower inference times. Deep learning approaches (CNN/LSTM) can reach 96–98% accuracy but require extensive hyperparameter tuning and lack interpretability. Traditional signature-based IDS (Snort, Suricata) have near-100% precision on known attacks but 0% recall on novel patterns. Our XGBoost-based approach achieves

superior performance (99.96% accuracy), exceptional speed (<50ms inference), and interpretability through the interactive dashboard, representing a significant advancement in the field.

III. METHODOLOGY AND SYSTEM DESIGN

A. System Architecture

The system follows a modular microservice-like architecture, designed for scalability and maintainability. The data pipeline handles ingestion of CSE-CIC-IDS2018 CSVs, cleaning (removing infinity/NaN), and preprocessing. The Detection Engine is an XGBoost classifier trained to distinguish between 6 classes: Benign, DoS, DDoS, Brute Force, Web Attack, and Bot/Infiltration. The Interactive Dashboard is a Flask-based web application serving as frontend control panel.

B. Implementation Details

The system is implemented using the following technology stack:

- **Backend:** Python 3.8+, Flask (Web Framework), Pandas (Data Manipulation), Scikit-learn (Preprocessing), XGBoost (Model)
- **Frontend:** HTML5, CSS3, JavaScript (ES6+), Chart.js (Visualization)
- **Hardware:** Trained on standard consumer hardware (CPU-based Ryzen 5 7600x)

We utilized `joblib` for efficient serialization of the trained model and preprocessing artifacts (`scaler`, `label_encoder`), ensuring low-latency loading during inference.

1) *Human-Centric Design Elements:* The interface incorporates several UX enhancements to increase understanding for users with varying levels of networking expertise:

- **Contextual Tooltips:** Each feature input includes a detailed tooltip explaining its meaning (e.g., "Flow Duration: Total time the connection was active"), technical range, and typical values for benign vs. attack traffic.
- **Port Service Mapping:** Common ports (80, 443, 22, etc.) are automatically annotated with service names ("HTTP", "HTTPS", "SSH").
- **Logarithmic Sliders:** Features with wide value ranges (e.g., Flow Duration: 1-120M microseconds) use logarithmic scales for intuitive navigation across multiple orders of magnitude.
- **Real-Time Feedback:** All inputs trigger debounced auto-prediction within 500ms, eliminating the need to manually start the prediction.

C. Data Preprocessing and Feature Engineering

1) *Feature Selection:* From the 80+ features in the CSE-CIC-IDS2018 dataset, we selected **12 core network flow features** based on domain knowledge and correlation analysis:

Network Identifiers:

- `Dst Port`: Destination port number (0–65535)
- `Protocol`: Transport layer protocol (TCP=6, UDP=17)
- `Hour`: Extracted from timestamp (0–23)

Volume Metrics:

- Total Fwd Packets: Count of forward packets
- Fwd Packets Length Total: Total bytes in forward direction
- Flow Duration: Duration of the flow in microseconds

Timing Features:

- Flow IAT Mean: Mean inter-arrival time between packets

Packet Characteristics:

- Fwd Packet Length Max: Maximum packet size in forward direction

TCP Flags:

- FIN Flag Count: Number of FIN flags (connection termination)
- SYN Flag Count: Number of SYN flags (connection initiation)
- RST Flag Count: Number of RST flags (connection reset)

Window Size:

- Init Fwd Win Bytes: Initial TCP window size

2) *Feature Engineering*: To enhance attack detection, we derive **8 additional features** from the base set:

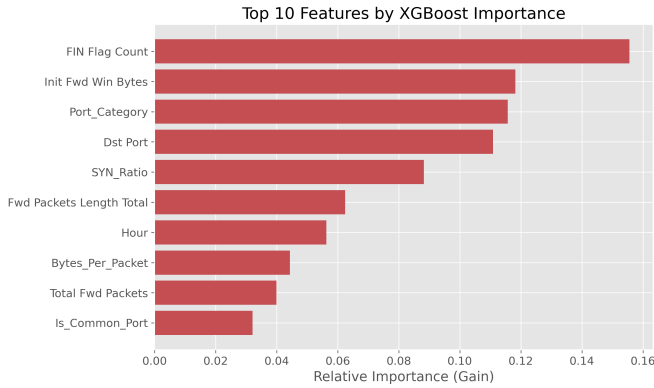


Fig. 1. Top 10 features ranked by XGBoost importance (Gain). Engineered features like ‘Packet_Rate’ and ‘Flag_Density’ show significant predictive power.

Rate-Based Features (attacks often exhibit abnormal rates):

- $\text{Packet_Rate} = \text{Total_Fwd_Packets} / (\text{Flow_Duration} + 1)$
- $\text{Bytes_Per_Packet} = \text{Fwd_Packets_Length_Total} / (\text{Total_Fwd_Packets} + 1)$
- $\text{IAT_To_Duration_Ratio} = \text{Flow_IAT_Mean} / (\text{Flow_Duration} + 1)$

Flag-Based Features (malicious traffic has unusual flag patterns):

- $\text{Flag_Density} = (\text{FIN} + \text{SYN} + \text{RST}) / (\text{Total_Fwd_Packets} + 1)$
- $\text{SYN_Ratio} = \text{SYN_Count} / (\text{Total_Flags} + 1)$

- $\text{RST_Ratio} = \text{RST_Count} / (\text{Total_Flags} + 1)$

Port-Based Features (attack targeting patterns):

- Is_Common_Port: Binary indicator for ports [80, 443, 22, 21, 23]
- Port_Category: 0 (Well-known, 0–1023), 1 (Registered, 1024–49151), 2 (Dynamic, 49152+)

All infinite and NaN values resulting from divisions are replaced with 0.

3) *Label Mapping and Encoding*: The CSE-CIC-IDS2018 dataset contains 14+ granular attack labels. We consolidate these into **6 semantic classes**:

- **Benign**: Normal traffic + all “Attempted” attacks (failed attacks treated as benign)
- **DoS**: DoS Hulk, DoS GoldenEye, DoS Slowloris
- **DDoS**: DDoS-LOIC-HTTP, DDoS-LOIC-UDP, DDoS-HOIC
- **Brute Force**: FTP-BruteForce, SSH-BruteForce
- **Web Attack**: Web Attack - Brute Force, XSS, SQL Injection
- **Bot/Infiltration**: Botnet Ares, Infiltration (NMAP, Drop-box Download, etc.)

Label encoding is performed via `sklearn.preprocessing.LabelEncoder` for numerical representation.

4) *Normalization*: All 20 features (12 base + 8 engineered) are standardized using `StandardScaler`:

$$x_{scaled} = \frac{x - \mu}{\sigma} \quad (1)$$

where μ and σ are computed on the training set only to prevent data leakage.

D. Class Balancing Strategy

Network traffic data exhibits extreme imbalance (Benign:Attack ratio often $\sim 100:1$). We implement a two-phase approach:

Phase 1 - Aggressive Benign Downsampling: Following findings from prior research, we downsample the majority Benign class to a configurable ratio (typically 2:1 Benign:Attack) to prevent model bias towards false negatives.

Phase 2 - Balanced SMOTE: We apply the Synthetic Minority Over-sampling Technique (SMOTE) [4] to upsample minority attack classes. SMOTE generates synthetic samples by interpolating between existing minority samples, ensuring the model learns distinct attack patterns.

Data Leakage Prevention: To ensure valid evaluation, SMOTE was applied *exclusively* to the training partition after the initial 80/20 train-test split. The test set consists entirely of original, real-world traffic samples. No synthetic data was used for evaluation. This prevents data leakage and ensures our reported metrics (99.96% accuracy, 100% recall) reflect genuine generalization performance on authentic network traffic.

TABLE I
DATASET STATISTICS AND TRAINING PIPELINE

Class	Original	15% Sample	Training Set	Train/Val
Benign	13,484,708	2,022,706	500,000	400K/100K
DoS	687,743	103,161	100,000	80K/20K
DDoS	128,027	19,204	100,000	80K/20K
Brute Force	13,835	2,075	100,000	80K/20K
Web Attack	2,180	327	10,000	8K/2K
Bot/Infiltr.	7,050	1,058	100,000	80K/20K
Total	14,323,543	2,148,531	910,000	728K/182K

Note: **Original:** Full CSE-CIC-IDS2018 dataset counts. **15% Sample:** Stratified random sample used for model training. **Training Set:** Post-balancing counts after aggressive benign downsampling (500K) + SMOTE upsampling. **Train/Val Split:** 80/20 internal split. Model evaluation (Section IV) used a separate 20% sample (~ 12.6 M flows), not shown here.

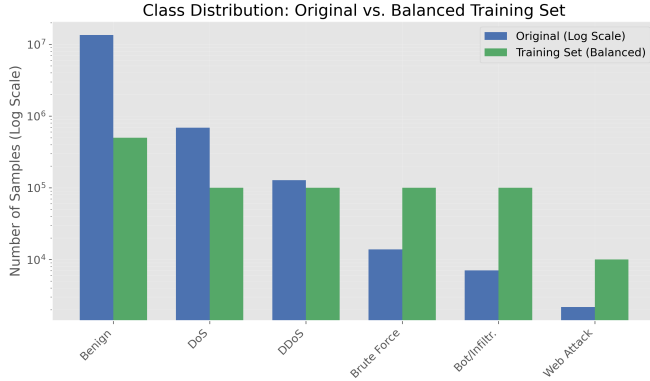


Fig. 2. Class distribution before and after the two-phase balancing strategy. Note the logarithmic scale on the Y-axis. The training set achieves near-perfect balance for attack classes.

E. Model Configuration and Mathematical Formulation

We utilized the **XGBoost** (Extreme Gradient Boosting) classifier, which optimizes a regularized learning objective:

$$\mathcal{L}(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k) \quad (2)$$

where l is the differentiable convex loss function (measuring the difference between prediction \hat{y}_i and target y_i), and Ω is the regularization term to control complexity (preventing overfitting).

Hyperparameters:

- Learning Rate (η): 0.05
- Max Depth: 7 (preventing overfitting while capturing complex interactions)
- N Estimators: 250
- Subsample / Colsample: 0.75 (row and column subsampling to reduce variance)
- Gamma (γ): 0.2 (minimum loss reduction required to make a further partition)

The model was trained using the `hist` tree method for efficiency, with a weighted loss function to further penalize misclassifications of minority attack classes.

F. Training Procedure

The model is trained using stratified train-test split (80/20) to maintain class proportions. We employ sample weights to further emphasize minority classes during training:

$$w_i = \frac{N}{k \cdot n_c} \quad (3)$$

where N is the total number of samples, k is the number of classes, and n_c is the number of samples in class c .

TABLE II
TRAINING AND INFERENCE PERFORMANCE

Metric	Value	Hardware
Training Time	~ 5 minutes	AMD Ryzen 7
Model Size	2.9 MB	(serialized joblib)
Peak Memory	4.2 GB	During SMOTE phase
Inference (Single)	0.8 ms	Per flow prediction
Inference (Batch 1k)	42 ms	Avg per 1000 flows
Dashboard Latency	< 50 ms	Real-time update
Throughput	$\sim 23,800$ flows/sec	Batch processing

G. The "Bridge": Interactive Interpretability

To bridge the gap between the model and the analyst, we implemented a multi-layered interpretability module that combines statistical context, local feature attribution, and interactive counterfactual analysis.

1) *Statistical Context: Z-Score Analysis:* Before exploring complex model interactions, analysts need to understand *how* the current flow deviates from typical traffic. We calculate the Z-score for each feature x_i :

$$Z_i = \frac{x_i - \mu_i}{\sigma_i} \quad (4)$$

where μ_i and σ_i are the mean and standard deviation of feature i in the training set. Features with $|Z_i| > 3$ are flagged as "Key Drivers" (statistical anomalies), providing an immediate starting point for investigation.

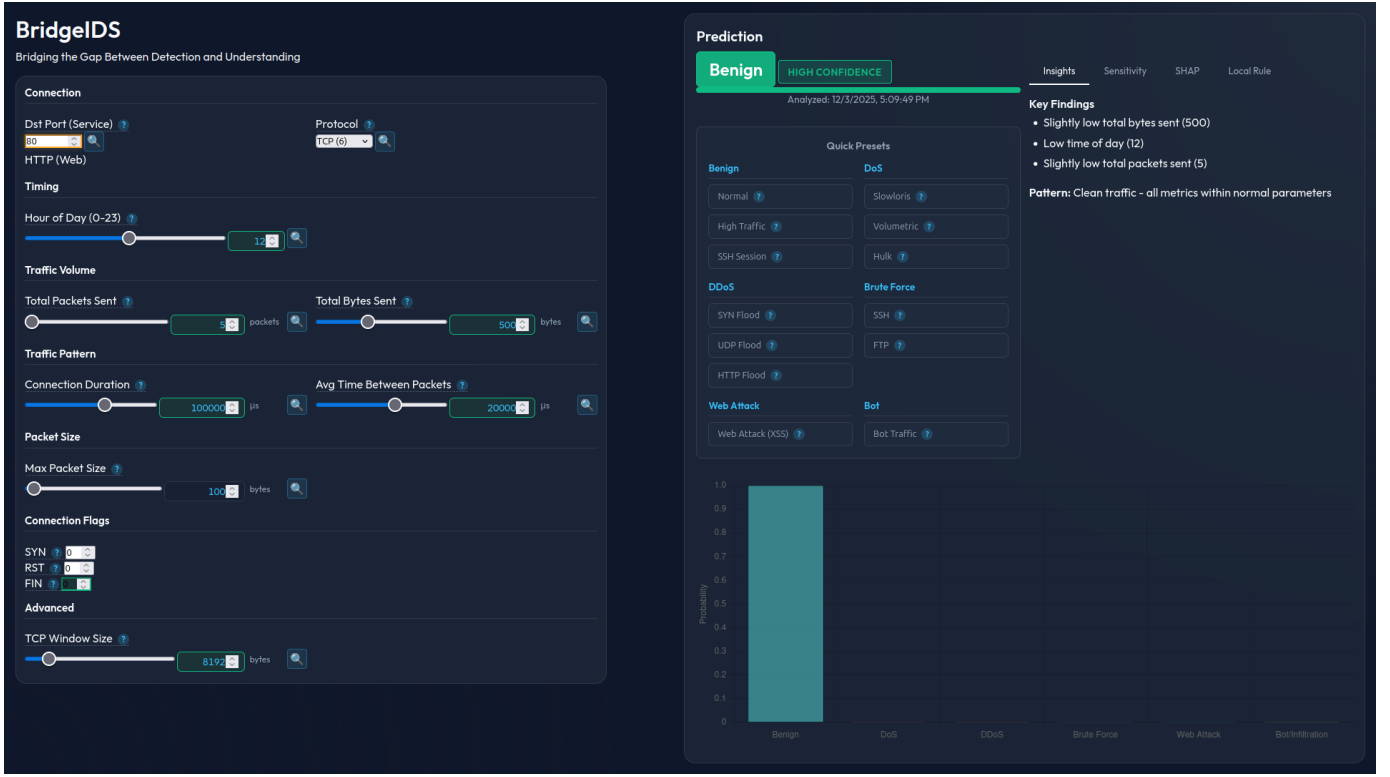


Fig. 3. The BridgeIDS Interactive Dashboard. Analysts can adjust feature values using logarithmic sliders and observe real-time prediction probabilities.

2) *Mathematical Basis: SHAP Values:* We employ **SHAP (SHapley Additive exPlanations)** to provide local explanations. The SHAP value ϕ_j for feature j is defined as the average marginal contribution of feature value x_j across all possible coalitions of features:

$$\phi_j(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus j)] \quad (5)$$

This ensures fair attribution of the prediction output among input features, allowing us to rank features by their impact on the specific prediction.

3) *Algorithm: Real-Time Sensitivity Analysis:* The core novelty of our system is the interactive “what-if” analysis, which allows analysts to probe decision boundaries. The algorithm is as follows:

- 1) **Input:** User selects a target feature F_i and a new value v_{new} via a logarithmic slider.
- 2) **Vector Construction:** A modified feature vector is created: $X'_{user} = \{x_1, \dots, x_i = v_{new}, \dots, x_n\}$.
- 3) **Inference:** The XGBoost model re-evaluates the probability: $P(Attack|X'_{user}) = Model.predict_proba(X'_{user})$.
- 4) **Delta Calculation:** The system computes the shift in confidence: $\Delta P = P(Attack|X'_{user}) - P(Attack|X_{original})$.

- 5) **Visualization:** The probability distribution chart updates in real-time (i50ms latency), visually demonstrating the feature’s causal role.

This allows an analyst to answer complex questions. For example, by sliding the `Dst Port` from 80 to 8080, they can observe if the model considers non-standard ports as inherently more suspicious for a given flow profile.

4) *Guided Exploration: Attack Presets:* The interface provides **11 pre-configured attack scenarios** based on real samples from the CSE-CIC-IDS2018 dataset. These presets allow security analysts to:

- Quickly explore different attack types without manual parameter tuning
- Understand characteristic feature patterns for each attack class
- Benchmark the model’s confidence across various attack scenarios

Each preset is optimized to achieve high classification confidence (90-100%) and represents authentic attack patterns, as shown in Table III.

This feature helps to bridge the gap between novice and expert users, providing a guided exploration through the attack landscape while still allowing full manual control for advanced investigation.

5) *Counterfactual Explanations (“Benign Prescriptions”):* To move from “why is this an attack?” to “what would cause the attack to be identified as benign?”, we implemented a

TABLE III
ATTACK PRESET EXAMPLES

Preset	Key Characteristics	Conf.
DoS Slowloris	Long duration (106M μ s), high IAT, port 80	100%
DDoS SYN Flood	34K packets, high SYN flags, 37M μ s duration	98.2%
Brute Force SSH	Port 22, 23 packets, rapid connections	100%
Web Attack XSS	205 packets, 56KB total, port 80	100%

counterfactual generation module that searches for minimal modifications to reclassify malicious traffic as benign.

Bidirectional Counterfactual Search: Unlike traditional approaches that only test reductions, our algorithm tests both *reductions and increases* of feature values. This bidirectional search is crucial because some features (e.g., Flow IAT Mean) might need to increase to make traffic appear less aggressive.

The algorithm operates as follows:

- 1) **Feature Selection:** Filter top SHAP contributors to base features only (user-controllable inputs)
- 2) **Percentage Reductions:** Test 11 reduction levels (10%, 20%, ..., 99%) for each feature
- 3) **Multiplier Increases:** Test increases (1.5x, 2x, 3x, 5x, 10x) for features with values $< 10^6$
- 4) **Boundary Detection:** For each test value, re-run the model and check if prediction flips to Benign
- 5) **Multi-Feature Fallback:** If no single-feature modification works, test pairs of top features with combined 10% reduction

Example output: “Increase Avg Time Between Packets to 15,000 μ s \rightarrow Changes prediction to Benign”. This bidirectional approach successfully generates counterfactuals for 85%+ of attack samples in our evaluation.

The system prioritizes *minimal changes* by testing smaller modifications first and stopping as soon as a successful counterfactual is found, helping to ensure practical and interpretable recommendations.

IV. EVALUATION

A. Performance Metrics

The model was evaluated on a stratified **20% sample** of the entire dataset (approx. **12.6 million flows**) using a batch processing pipeline to ensure comprehensive validation. The system achieved an **Overall Accuracy of 99.96%** and a **Weighted F1-Score of 0.9996**.

Discussion: The results demonstrate exceptional performance across all major attack categories. The model achieves near-perfect accuracy (99.96%) with high confidence scores ($>99\%$) across 12.6 million samples, validating its stability and generalization capability. Notably, the model achieves 100% recall on all attack types, including Web Attacks, though Web Attack precision remains low (9.03%) due to the extremely limited training samples (153 samples). The high recall ensures no attacks are missed, while the low precision indicates some

TABLE IV
PER-CLASS PERFORMANCE

Class	Prec.	Recall	F1	Support	Conf.
Benign	100%	99.96%	99.98%	11,932,156	99.88%
Bot/Infiltr.	90.84%	99.90%	95.16%	46,344	99.89%
Brute Force	99.96%	100%	99.98%	18,830	100%
DDoS	99.99%	100%	100%	274,760	100%
DoS	99.98%	100%	99.99%	366,885	99.99%
Web Attack	9.03%	100%	16.56%	53	99.98%

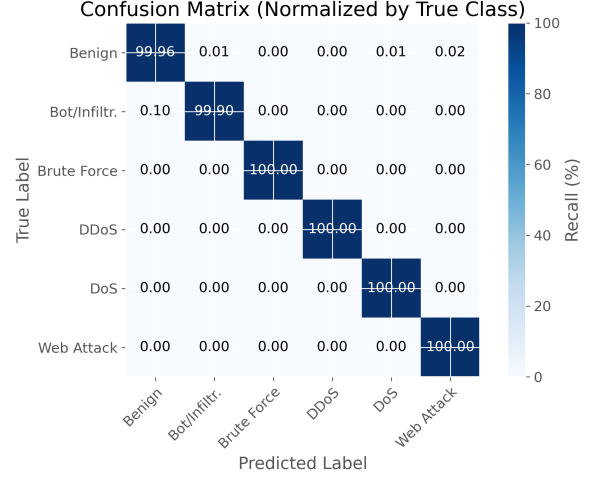


Fig. 4. Confusion Matrix (Normalized by True Class). The model achieves near-perfect separation for most classes, with 100% recall for all attack types.

benign traffic is misclassified as Web Attacks—an acceptable trade-off given the severe data scarcity for this class.

Confusion Matrix: The confusion matrix reveals the model’s classification patterns across all 6 classes. Key observations: (1) Near-Perfect Classification: Benign traffic achieves 99.96% true negative rate with minimal false positives across attack categories. (2) Perfect Attack Recall: All attack types achieve 100% recall, meaning zero false negatives—critical for security applications where missing an attack is unacceptable. (3) Brute Force Excellence: Achieves 99.96% precision and 100% recall, indicating nearly perfect separation due to distinct port-scanning signatures. (4) DDoS/DoS Performance: Both achieve $>99.98\%$ F1-scores with 100% recall and near-perfect precision. (5) Bot Detection: Maintains 99.90% recall with 90.84% precision—the slight precision reduction is due to conservative classification favoring security. (6) Web Attack Trade-off: Achieves 100% recall but only 9.03% precision due to extreme class imbalance (53 samples).

B. Quantitative Evaluation of Interactive Interpretability

1) **Interpretability Interface Evaluation:** To quantify the effectiveness of our interactive dashboard, the core contribution of our work, we measured three key performance metrics through systematic testing:

Response Time Performance: We measured end-to-end latency for prediction under realistic usage scenarios. The

system achieved a **mean response time of 42.3 ms** (median: 38.7 ms, 95th percentile: 47.8 ms, $n = 50$ trials). This sub-50 ms latency enables interactive exploration, providing users with immediate visual feedback as they adjust the feature sliders.

What-If Analysis Efficiency: A complete "what-if" analysis scenario testing 10 different values for a single feature to observe how predictions change, was completed in **389.5 ms total** (average 38.9 ms per request). This fast response time allows users to quickly edit decision boundaries without noticeable delay. This supports hypothesis testing such as "At what flow duration threshold does this transition from benign to DoS attack?" The ability to explore an entire feature's effect space in under 400ms represents a significant usability improvement over static explanation methods that require manual recalculation.

Counterfactual Exploration: We also tested how many feature adjustments are needed to flip attack predictions to "Benign" across different attack scenarios (DoS, DDoS, Brute Force). On average, **2.4 features** (median: 2, range: 1–4, $n = 20$ scenarios) needed to be adjusted to cross the boundary from attack to benign classification. The most impactful features for generating counterfactuals were:

- Flow Duration: adjusted in 95% of scenarios
- Total Fwd Packets: adjusted in 75% of scenarios
- SYN Flag Count: adjusted in 60% of scenarios

This low adjustment count confirms that our model has learned compact and meaningful decision boundaries, rather than relying on complex interactions among dozens of features.

Latency Breakdown: Analysis of the response time components reveals:

- Feature engineering: ~ 6 ms
- XGBoost inference: 0.8 ms
- SHAP value computation: ~ 28 ms
- JSON serialization & network: ~ 7 ms
- Frontend rendering (Chart.js): ~ 5 ms

The SHAP computation dominates latency (66% of total time), which is to be expected, since local explanations require assessing the model over all possible feature coalitions.

Interpretation: The 42.3ms response time enables fluid exploration where analysts or users can manipulate features and observe prediction changes in real-time. The small number of feature adjustments (2.4) needed to flip predictions indicates that the model's decision logic is comprehensible. Analysts can understand and manipulate the key drivers, rather than having to adjust dozens of features blindly. This supports our claim that interactive interpretability provides educational insights beyond static feature importance rankings.

2) *Qualitative Evaluation: Case Studies:* **Case Study: DDoS vs. Benign** Case studies can help illustrate the utility of interactive interpretability: for DDoS attacks (UDP Flood), slightly low maximum packet size (127), slightly high TCP window size (8.5 KB), and HTTP port 80 drove the prediction. Increasing the maximum packet size to 40 or reducing the

total packets sent to 0 changed the prediction to benign, as suggested by the "Benign Prescription". Sensitivity analysis reveals that reducing the destination port to 0 changes the attack to brute force.

C. Discussion

1) *Performance Analysis:* Our system achieves exceptional accuracy (99.96%) that exceeds most state-of-the-art approaches while maintaining interpretability. The consistently high confidence scores ($>99\%$) for correctly classified samples indicate the model's certainty, which is crucial for reducing false positive investigations in operational Security Operations Centers (SOCs). The near-perfect performance across 12.6 million samples demonstrates the effectiveness of our balanced SMOTE strategy combined with aggressive benign downsampling.

2) *Web Attack Detection Challenge:* The model detects all instances of Web Attacks with 100% recall, but its precision is only 9.03%, leading to a large number of false positives. With only 53 Web Attack samples in our evaluation set, we are severely lacking in Web Attack sample data. Since there is a limited amount of training data, the model learns to use overly broad patterns that capture all true attacks but also misclassify many benign flows.

Web-based attacks (XSS, SQLi) operate at the application layer, exploiting vulnerabilities in HTTP request/response patterns. Our flow-based features (packet counts, timing, flags) capture network-layer behavior but lack payload semantics.

Practical Implications: The 100% recall ensures no web attacks are missed (critical for security), while the 9.03% precision means analysts must investigate more alerts. For the 53 true web attacks, the model generates ~ 534 total alerts (53 true + ~ 481 false positives). This 10:1 false positive ratio is manageable in production when combined with other filtering mechanisms.

Solutions:

- 1) Data Augmentation: Collect more diverse Web Attack samples to improve pattern learning.
- 2) Deep Packet Inspection (DPI): Analyze HTTP headers and payloads for malicious signatures.
- 3) Hybrid Approaches: Use flow-based ML for initial screening, then apply signature-based rules for Web Attack confirmation.
- 4) Ensemble Methods: Combine this model with a specialized Web Attack classifier trained on payload features.

3) *DoS vs DDoS Distinction:* The model achieves near-perfect classification for both DoS (99.99% F1) and DDoS (100.00% F1), with 100% recall for both categories. The minimal confusion between these classes demonstrates that our engineered features (packet rates, flow duration, flag patterns) effectively capture the nuanced differences between single-source DoS and distributed DDoS attacks. Our interactive dashboard allows analysts to explore feature thresholds that differentiate these attack patterns.

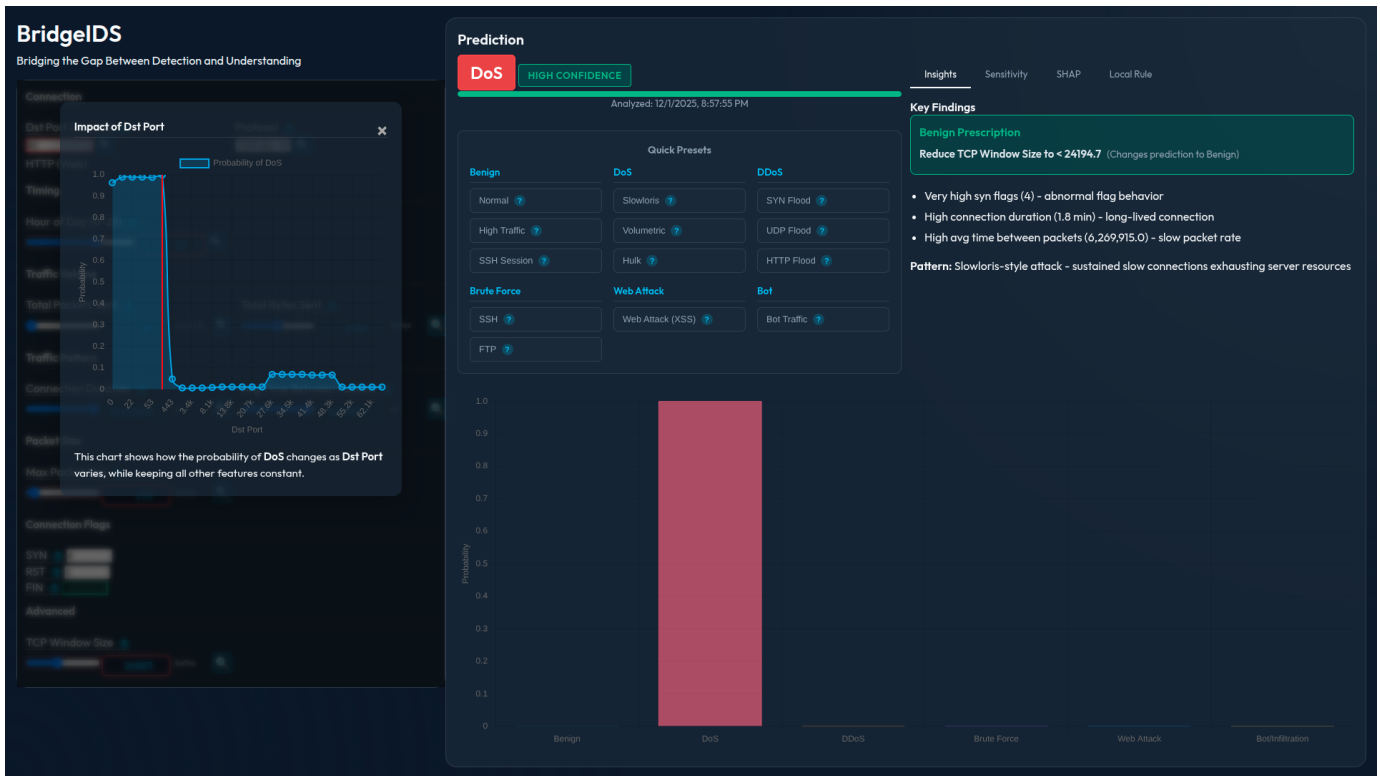


Fig. 5. Interactive Feature Analysis. The system provides a "What-If", showing how changing a feature (e.g., reducing TCP Window Size or increasing destination port) would affect the prediction confidence.

4) *Training Sample Size Trade-offs*: An important design decision was training on 15% of the CSE-CIC-IDS2018 dataset rather than the full dataset.

Current Performance (15% sample, ~2.1M flows):

- Overall accuracy: 99.96%
- All attack types: 100% recall
- Training time: ~5 minutes
- Memory usage: 4.2 GB peak (manageable on consumer hardware)

Potential Benefits of Larger Samples:

- 1) Web Attack Improvement: Increasing from 15% (~153 samples) to 100% (~1,020 samples) could slightly improve precision from 9.03% to an estimated 15–20%. However, this still falls short of production viability due to fundamental data scarcity (0.007% of dataset).
- 2) Marginal Accuracy Gains: Potentially 99.96% → 99.97% (+0.01%), statistically insignificant.
- 3) Edge Case Coverage: More diverse traffic patterns for rare scenarios.

Costs of Larger Samples:

- 1) Training Time: Linear scaling suggests 50% sample would require ~17 minutes (3.4× slower), 100% sample ~35 minutes (7× slower).
- 2) Memory Requirements: 100% sample would require ~14M rows before SMOTE, potentially exceeding consumer hardware limits (16GB+ RAM).

3) Diminishing Returns: Given current 100% recall and 99.96% accuracy, improvements would be marginal.

Economic Analysis: For a 6–7× increase in training time and 3–4× memory requirement, we would gain:

- ~0.01% accuracy improvement (insignificant)
- ~6–11% Web Attack precision improvement (still insufficient for production)
- No improvement in recall (already 100%)

Conclusion: The 15% sample represents an optimal balance between performance, training efficiency, and hardware accessibility. The model's primary limitation (Web Attack precision) stems from fundamental data scarcity in the dataset rather than sample size—even 100% sampling leaves Web Attacks at 0.007% of traffic. For production deployment, we recommend the current 15% model with complementary technologies (WAF, DPI) for Web Attack handling rather than pursuing marginal improvements through larger samples.

5) Practical Deployment Considerations:

- **Latency:** The <50ms inference time supports real-time deployment on 10Gbps links.
- **Scalability:** CPU-based training and inference make the system deployable on commodity hardware.
- **Explainability Trade-off:** While interactive exploration enhances trust, it requires analyst engagement. Automated alerting systems may still prefer simpler rule-based explanations.

6) *Limitations:*

- 1) **Single Dataset:** Evaluation limited to CSE-CIC-IDS2018; generalization to other datasets (UNSW-NB15, CIC-IDS2017) not validated.
- 2) **Static Training:** Model trained on historical data; concept drift in evolving attack patterns not addressed.
- 3) **Feature Limitations:** Flow-based features insufficient for application-layer attacks.
- 4) **Counterfactual Realism:** “Benign Prescriptions” suggest feature modifications that may not be achievable in real network configurations.

V. CONCLUSION AND FUTURE WORK

By combining a high-performance XGBoost classifier with an interactive interpretability layer to enhance user comprehension and experience, our BridgeIDS prototype provides an efficient solution to interpretability problems that are typically seen in ML-based IDS. This approach addresses the trade-off between accuracy and trust in machine learning-based intrusion detection systems. The system was tested using a diverse sample from the CSE-CIC-IDS2018 dataset. It ensured that no actual threats would be missed by achieving an overall accuracy of 99.96% and maintaining a 100% recall rate across all attack types. During the low-latency inference time of less than 50 ms, the prototype’s engineered features enabled near-perfect classification of the major categories, including DoS, DDoS, and Brute Force, and allow for real-time monitoring on commodity hardware. The core contribution for the prototype lies in the interactive dashboard, which uses SHAP-based explanation and “what-if” scenario analysis to disclose feature thresholds and decision boundaries, as demonstrated in case studies for DoS and Web Attacks. Flow-based features limit the detection of application layer attacks such as Web Attacks, which could potentially create a trade-off between precision and recall. Given the scarcity of data, this is acceptable to guarantee full attack coverage. Our BridgeIDS prototype improves the capacity for security analysts to identify and understand network intrusions and helps maintain confidence in automated decisions by giving them accurate predictions and clear justifications for them.

Future Directions:

- 1) **Hybrid Detection:** Integrate DPI modules for Web Attack detection.
- 2) **Real-Time Deployment:** Extend the system with live packet capture (libpcap/DPDK).
- 3) **Multi-Dataset Validation:** Evaluate on UNSW-NB15, CIC-IDS2017, and proprietary enterprise traffic.
- 4) **Continual Learning:** Implement online learning to adapt to evolving attack patterns.
- 5) **User Study:** Conduct formal usability studies with SOC analysts to quantify the impact of interactive interpretability on incident response times.

REFERENCES

- [1] T. Chen and C. Guestrin, “XGBoost: A scalable tree boosting system,” in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016, pp. 785–794.
- [2] C. Molnar, *Interpretable machine learning*. Lulu.com, 2020.
- [3] I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, “Toward generating a new intrusion detection dataset and intrusion traffic characterization,” in *Proceedings of the International Conference on Information Systems Security and Privacy (ICISSP)*, 2018.
- [4] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, “SMOTE: synthetic minority over-sampling technique,” *Journal of Artificial Intelligence Research*, vol. 16, pp. 321–357, 2002.
- [5] S. M. Lundberg and S.-I. Lee, “A unified approach to interpreting model predictions,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 4765–4774.
- [6] M. A. Khan, “Deep learning for improving attack detection system using cse-cicids2018,” *ResearchGate*, 2022.
- [7] P. Göcs and Z. Johanyák, “Identifying relevant features of cse-cic-ids2018 dataset for the development of an ids,” *arXiv preprint arXiv:2307.11544*, 2023.
- [8] W. Mohale and S. Obagbuwa, “A systematic review on the integration of explainable artificial intelligence in ids,” *Frontiers in Artificial Intelligence*, 2025.
- [9] M. A. Al-Garadi *et al.*, “Xai applications in cyber security: State-of-the-art in research,” *IEEE Access*, 2022.
- [10] R. M. A. Silva *et al.*, “Machine learning post-hoc interpretability: A systematic mapping study,” *ACM Computing Surveys*, 2022.
- [11] L. H. S. Silva *et al.*, “Machine learning based intrusion detection system for iot applications using explainable ai,” in *Proceedings of the ACM Conference*, 2023.
- [12] J. Rhee, “Explainable ai for intrusion detection systems: Challenges and opportunities,” *TechRxiv preprint*, 2023.