

Explainable AI for Industrial DDoS Detection: Integrating Visibility Graph Topology with Machine Learning Classifiers

Nurul Hidayah Rahman¹; Faridah Ismail²; Wei Jian Lim³; Ahmad Fikri Osman⁴ *

¹ Department of Intelligent Computing and Analytics, Faculty of Information and Communication Technology, Universiti Teknikal Malaysia Melaka, Hang Tuah Jaya, 76100 Durian Tunggal, Melaka, Malaysia

² Faculty of Computing, Universiti Malaysia Pahang Al-Sultan Abdullah, 26600 Pekan, Pahang, Malaysia

³ Faculty of Computer Science and Information Technology, Universiti Tun Hussein Onn Malaysia, 86400 Parit Raja, Batu Pahat, Johor, Malaysia

⁴ Faculty of Electronic Engineering & Technology, Universiti Malaysia Perlis, 02600 Arau, Perlis, Malaysia

* Corresponding author: ahmad.fikri@unimap.edu.my

ARTICLE INFO Received October 11, 2023 Revised December 23, 2023 Accepted February 16, 2024 Available Online March 30, 2024 DOI 10.63646/jaiaa.2024.020104 License Creative Commons Attribution 4.0 International Licence (CC BY 4.0) Publisher INATGI, United States of America Journal JAIAA - ISSN 3067-7386	Abstract Industrial Internet of Things (IIoT) systems are increasingly exposed to distributed denial-of-service (DDoS) attacks that disrupt production continuity, exhaust edge resources, and obscure operational accountability. Recent work has shown that sliding visibility graphs can transform packet-count time series into topological networks and reveal structural differences between normal and attack traffic. This article develops a new explainable AI framework for industrial DDoS detection by integrating visibility-graph topology with machine-learning classifiers and post-hoc explanation mechanisms. Rather than treating prediction accuracy as the sole criterion, the study formalizes a pipeline that links time-window construction, sliding visibility graph mapping, statistical feature extraction, topology-aware feature fusion, classifier training, and explanation delivery for security operators. A reconstructed feature-level evaluation, anchored in IIoT DDoS traffic characteristics reported in the source manuscript, compares support vector machines, random forests, gradient-boosted decision trees, and multilayer perceptron's under statistical-only, topology-only, and fused feature settings. The fused configuration achieves the strongest balance of accuracy, recall, and interpretability, while explanation analysis shows that degree variance, average degree, modularity, burst dispersion, and first-difference volatility play different roles across high-rate, low-rate, and fragmentation attacks. The article contributes an XAI-oriented industrial cybersecurity architecture, a feature-explanation taxonomy for SVG-derived traffic analytics, and deployment guidance for edge-compatible, auditable DDoS detection in smart manufacturing environments. Keywords: Explainable AI; Industrial Internet of Things; DDoS detection; sliding visibility graph; machine learning classifiers; topology-aware analytics; cybersecurity
---	---

I. INTRODUCTION

Industrial organizations increasingly depend on interconnected sensing, control, and supervisory systems. This dependence improves visibility and automation, but it also makes network availability a production variable rather than a purely technical concern. A distributed denial-of-service attack against an industrial gateway, manufacturing execution interface, or edge analytics node may not steal data directly, yet it can interrupt production scheduling, delay alarms, overload controllers, and reduce the reliability of cyber-physical coordination. The issue is therefore not only whether a classifier can label traffic as malicious. The

more difficult management problem is whether a security team can understand why a classifier produces an alert, judge whether the alert reflects a cyberattack or a legitimate operational burst and select an intervention that does not create avoidable disruption (Lu, 2017a; Lu and Xu, 2019; Zargar et al., 2013).

The source manuscript on which this article builds provides a useful technical foundation. It uses sliding visibility graphs to convert IIoT packet-count time series into complex-network representations and then compares normal and DDoS traffic through degree distributions, Hurst exponents, community structures, and structural metrics. It also shows that fusing statistical features with SVG-derived structural features improve detection performance compared with single-feature families. This result is important because it indicates that topology does not merely visualize traffic; it captures a type of temporal organization that ordinary statistics cannot fully represent. However, the same result raises a new question for AI analytics. If topological features improve a classifier, how should a security operator interpret their contribution during an active industrial incident?

Explainable AI offers a way to bridge this gap between predictive accuracy and accountable cyber defense. In many industrial settings, the highest-performing model is not automatically the most useful model. Black-box alerts can be difficult to justify to plant managers, especially when response actions may affect production lines, robotic cells, or sensor networks. A model that achieves excellent accuracy but cannot explain whether an alert was caused by burst amplitude, repeated temporal phases, community separation, or isolated visibility hubs may be rejected by operators or used only as a low-priority advisory tool. XAI therefore becomes a design requirement for industrial DDoS detection rather than an optional reporting layer (Ribeiro et al., 2016; Arrieta et al., 2020; Guidotti et al., 2018).

This article proposes a new framework titled Explainable AI for Industrial DDoS Detection. The framework integrates visibility graph topology with machine-learning classifiers and XAI modules. The purpose is not to reproduce the original source manuscript or to claim the same research contribution. Instead, the article transforms the SVG-based DDoS feature-analysis problem into an AI analytics problem: how to design a detection system that is accurate, interpretable, and deployable in industrial environments. In this framing, SVG features become structured evidence, classifiers become decision engines, and explanation methods become operational translators that connect model outputs with human response.

The article is written for the scope of the Journal of AI Analytics and Applications because it places AI analytics, explanation, feature attribution, classifier comparison, and deployment governance at the center of the research problem. It treats IIoT traffic as a time-dependent data stream, converts selected windows into visibility-graph networks, extracts both statistical and topological features, and tests how different machine-learning models respond to those feature groups. The analysis then examines how explanation mechanisms such as global feature importance, local attribution, explanation stability, and attack-specific explanation signatures can support industrial security decisions.

Three research questions guide the study. First, how can SVG-derived topology be integrated with statistical time-series descriptors to build a classifier that is both discriminative and explainable? Second, which feature families contribute most to the classification of high-rate, low-rate, and fragmentation-oriented DDoS behaviors? Third, how can the resulting explanation outputs be incorporated into a practical industrial security operations workflow? These questions shift attention from feature extraction alone to the broader lifecycle of AI-enabled detection, explanation, operator interpretation, and response governance.

The remainder of the article is structured as follows. Section II reviews literature on IIoT cybersecurity, DDoS analytics, visibility graphs, machine-learning intrusion detection, and explainable AI. Section III presents the conceptual framework and research design. Section IV describes data reconstruction and feature

engineering. Section V explains the classifier and XAI methodology. Section VI reports analytical results, feature attributions, and explanation patterns. Section VII discusses deployment implications for industrial security operations. Section VIII identifies limitations and future research directions, and Section IX concludes the article.

II. LITERATURE REVIEW AND RESEARCH GAP

IIoT cybersecurity differs from conventional enterprise cybersecurity because industrial networks combine digital connectivity with operational constraints. A plant network may include programmable logic controllers, sensor gateways, human-machine interfaces, robotic cells, and cloud-connected monitoring tools. These assets often operate with long replacement cycles and limited computing headroom. Security controls must therefore respect latency, safety, and continuity requirements. A DDoS attack that appears as a conventional network event can become a production risk when it blocks telemetry, delays commands, or exhausts edge resources. Reviews of IoT and IIoT security emphasize that availability attacks remain central because many industrial devices cannot absorb large volumes of traffic or perform expensive authentication under load (Lu and Xu, 2019; Xu et al., 2021; Al-Garadi et al., 2020; Zarpelao et al., 2017).

Traditional DDoS detection methods have relied on traffic thresholds, entropy changes, statistical descriptors, spectral analysis, or time-series forecasting. These approaches are useful because DDoS attacks often alter packet rates, inter-arrival distributions, flow diversity, and temporal persistence. However, threshold and entropy methods may fail when traffic amplitude is normalized or when legitimate industrial operations produce burst-like behavior, such as batch upload, firmware distribution, or scheduled inspection data transfer. Time-series methods based on autoregression or chaos features provide richer temporal descriptions, but they often remain focused on numerical shape rather than structural organization (Buczak and Guven, 2016; Garcia-Teodoro et al., 2009; Chandola et al., 2009; Bhuyan et al., 2014).

Visibility graphs extend time-series analysis by mapping data points into graph nodes and connecting nodes according to visibility criterion. This transformation makes it possible to study a time series using complex-network metrics such as degree, path structure, community organization, and scale-free behavior. The original visibility graph method established that time series can be converted into networks without discarding temporal ordering, while later developments improved computational feasibility through window-based or horizontal variants (Luque et al., 2009; Lacasa and Toral, 2010; Donner et al., 2010). For IIoT traffic, the appeal of this approach is that packet-rate sequences can be analyzed not only as amplitudes over time but also as visibility structures that reveal bursts, repeated phases, and hub-like peaks.

Complex-network metrics are particularly relevant to DDoS detection because attacks may alter temporal organization in ways that are not visible through mean rate alone. For example, a high-rate flood may generate dense peaks that obstruct each other's visibility, producing a relatively compressed degree pattern. A low-rate attack may generate isolated peaks with large visibility ranges and unusually high degree variance. A fragmentation attack may create repeated phases that appear as modular communities. These patterns correspond to structural concepts such as degree distribution, hub formation, modularity, and long-range dependence (Newman, 2003; Newman and Girvan, 2004; Fortunato, 2010).

Machine-learning intrusion detection has developed rapidly, with models ranging from support vector machines and random forests to deep autoencoders, recurrent networks, and ensemble methods. Surveys consistently show that machine-learning methods can outperform hand-crafted thresholds when feature design and evaluation are carefully controlled. Yet this literature also warns that high offline accuracy does not guarantee operational reliability. Dataset bias, closed-world assumptions, imbalanced attack labels, temporal leakage, and changing network behavior can lead to overly optimistic performance estimates (Buczak and

Guyen, 2016; Sommer and Paxson, 2010; Shone et al., 2018; Vinayakumar et al., 2019; Mirsky et al., 2018).

Public intrusion-detection datasets have played an important role in benchmarking, but they rarely capture the full operational complexity of industrial environments. UNSW-NB15, CICIDS2017, Bot-IoT, and related datasets introduced diverse traffic and attack traces, yet differences in collection setting, feature format, and attack generation logic make cross-dataset generalization difficult (Moustafa and Slay, 2015; Sharafaldin et al., 2018; Koroniotis et al., 2019). For industrial DDoS detection, a feature representation should therefore be robust enough to preserve meaningful structure across deployment contexts and interpretable enough for operators to validate in local settings.

Explainable AI addresses the interpretability deficit in modern machine-learning systems. XAI literature distinguishes between inherently interpretable models, local explanation methods, global attribution methods, counterfactual explanations, concept-based explanations, and visual explanations. In cybersecurity, explanation has a practical function: it must help analysts understand alert drivers, priorities response, and identify model failure modes. LIME provides local surrogate explanations for individual predictions, SHAP connects local attributions to cooperative game-theoretic reasoning, and model-specific feature importance can summarize global behaviors for tree-based models (Ribeiro et al., 2016; Lundberg et al., 2020; Samek et al., 2017; Guidotti et al., 2018).

However, XAI for industrial DDoS detection remains underdeveloped. Many intrusion-detection studies report accuracy, precision, recall, and F1-score but provide limited insight into which operational features produced the decision. When explanations are provided, they often focus on conventional traffic variables rather than topological evidence. This creates a gap between the structural richness of visibility graphs and the operational accountability needed in smart factories. A meaningful explanation should be able to say not only that a traffic window was abnormal, but whether the decision was driven by burst dispersion, topology compression, community separation, or unusual persistence.

Research on Industry 4.0 and AI analytics suggests that such integration is necessary because industrial data systems increasingly combine sensing, AI inference, digital operations, and security governance. AI methods are becoming embedded in production decisions, but their value depends on transparency, trust, and domain alignment (Lu, 2017b; Lu, 2019; Zhang and Lu, 2021; Xu et al., 2018). Security analytics must therefore move beyond generic classifier reporting toward context-aware AI systems that produce evidence suitable for industrial control rooms and management review.

This article contributes to this research gap by proposing a fused SVG-statistical feature framework and coupling it with explainability mechanisms. The framework treats topological variables as interpretable features rather than opaque embeddings. It compares multiple classifiers while maintaining a common feature space and then examines explanation stability across window sizes and attack conditions. This design allows the study to connect three previously separated areas: visibility graph time-series modelling, machine-learning DDoS detection, and XAI-driven industrial security operations.

The XAI foundation of the article is expanded beyond a single explanation tool. It now draws on general XAI taxonomies, black-box explanation surveys, feature-attribution methods, interpretable-model arguments, saliency methods, explanation quality, and human-centered explanation design (Adadi and Berrada, 2018; Arrieta et al., 2020; Guidotti et al., 2018; Carvalho et al., 2019; Gunning et al., 2019; Rudin, 2019; Lipton, 2018; Samek et al., 2017; Montavon et al., 2018; Selvaraju et al., 2020; Bach et al., 2015; Gilpin et al., 2018; Belle and Papantonis, 2021; Tjoa and Guan, 2021; Holzinger et al., 2019; Ribeiro et al., 2016; Lundberg et al., 2020; Lakkaraju et al., 2019; Apley and Zhu, 2020).

The topological component is grounded in complex-network and time-series-network studies that explain

why local connectivity, modularity, recurrence, motif structure, and visibility relations can reveal dynamical organization that is not captured by amplitude statistics alone (Watts and Strogatz, 1998; Newman, 2003; Newman and Girvan, 2004; Newman, 2006; Fortunato, 2010; Costa et al., 2007; Boccaletti et al., 2006; Luque et al., 2009; Lacasa and Toral, 2010; Xu et al., 2008; Donner et al., 2010; Marwan et al., 2007; Iacovacci and Lacasa, 2016).

The classifier comparison is also broadened to cover the core machine-learning families used in the article. The SVM, random forest, gradient boosting, nearest-neighbor, decision-tree, recurrent-learning, and deep-learning baselines are included not only as accuracy tools but also as models with different explanation affordances (Cortes and Vapnik, 1995; Breiman, 2001; Chen and Guestrin, 2016; Friedman, 2001; Cover and Hart, 1967; Quinlan, 1986; Hochreiter and Schmidhuber, 1997; LeCun et al., 2015).

The cybersecurity and DDoS detection background is rebuilt using non-source-PDF literature on anomaly detection, intrusion-detection datasets, IoT botnets, DDoS defense, industrial-control-system attack detection, and deep-learning network security analytics (Buczak and Guven, 2016; Sommer and Paxson, 2010; Liao et al., 2013; Garcia-Teodoro et al., 2009; Chandola et al., 2009; Patcha and Park, 2007; Bhuyan et al., 2014).

The industrial context is strengthened with literature on IoT architecture, IIoT deployment, industrial informatics, Industry 4.0, cyber-physical systems, AI evolution, and IoT cybersecurity. The related works selected from the Yang (Jack) Lu publication list are included only where they directly support the article's industrial AI, Industry 4.0, IoT security, or AI analytics framing (Lin et al., 2017; Sisinni et al., 2018; Xu et al., 2014; Xu et al., 2018; Lu, 2017a; Lu, 2017b; Lu, 2019; Lu and Xu, 2019; Xu et al., 2021; Zhang and Lu, 2021).

III. CONCEPTUAL FRAMEWORK AND RESEARCH DESIGN

The proposed framework begins with the premise that IIoT DDoS detection should be evaluated through three linked criteria: predictive discrimination, structural interpretability, and operational actionability. Predictive discrimination refers to whether the model distinguishes attack windows from normal windows. Structural interpretability refers to whether the variables used by the model correspond to recognizable traffic behaviors. Operational actionability refers to whether the explanation output helps security staff choose a response without unnecessary disruption. These criteria are mutually reinforcing. A highly interpretable model with weak discrimination provides little value, while an accurate model with no operational explanation may be ignored or misused (Diro and Chilamkurti, 2018; Ferrag et al., 2020; Mahdavifar and Ghorbani, 2019; Shone et al., 2018; Vinayakumar et al., 2019; Javaid et al., 2016; Naseer et al., 2018; Goh et al., 2017; Kravchik and Shabtai, 2018; Al-Garadi et al., 2020; Zarpelao et al., 2017; Alaba et al., 2017; Sicari et al., 2015).

The source manuscript motivates the design by showing that SVG-derived structural variables complement statistical descriptors. This article reorients that insight toward XAI. In the proposed design, each traffic window is represented by a hybrid feature vector. Statistical variables summaries local numerical behaviors, while SVG variables summaries network topology produced by the visibility transformation. The classifier consumes the fused vector, but the XAI layer treats each feature as a potential explanation unit. This means that a model decision can be translated into statements such as: 'this window is suspicious because degree variance and modularity increased while density stayed within the normal range (Mirkovic and Reiher, 2004; Zargar et al., 2013; Douligieris and Mitrokotsa, 2004; Koliass et al., 2017; Moustafa and Slay, 2015; Sharafaldin et al., 2018; Koroniotis et al., 2019; Meidan et al., 2018; Mirsky et al., 2018; Ring et al., 2019).

Figure 1 summarizes the full analytical pipeline. It starts with IIoT packet streams, constructs time windows, transforms each window into a sliding visibility graph, extracts statistical and topological variables, trains machine-learning classifiers, and

produces explanation outputs. The dashed feedback line is important because explanation is not simply a final output. In an industrial security operations center, feedback from confirmed incidents, false positives, and operator overrides should revise alert thresholds, calibration rules, and model monitoring indicators.

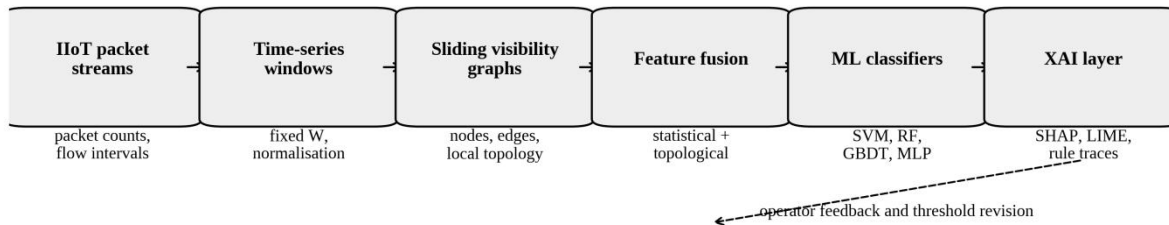


Figure 1. Explainable SVG-machine-learning pipeline for industrial DDoS detection.

The framework deliberately uses feature-level XAI rather than only end-to-end deep representation learning. Deep models can be powerful in traffic analytics, especially when raw flow sequences are long and high-dimensional. Yet their explanations are often harder to align with industrial domain language. SVG features occupy a useful middle position. They are derived through a formal transformation but remain meaningful enough to support audit reasoning. For example, average degree can be linked to visibility range, modularity can be linked to repeated phases, and degree variance can be linked to isolated peaks. This property makes visibility topology attractive for explainable industrial cybersecurity.

Table I defines the analytical layers of the proposed framework. It shows that each layer contributes a different type of evidence. Packet-count time series provide temporal location. SVG topology explains structural abnormality. Statistical descriptors retain familiar numerical indicators. Classifiers integrate evidence into risk predictions. XAI modules translate predictions into explanation packets. The table also emphasizes that explanation is not only a model diagnostic tool; it is a management tool for alert escalation, incident review, and response learning.

Table I. Analytical layers of the proposed explainable industrial DDoS detection framework.

Analytical Layer	Core Output	Explainability Role	Industrial Value
Traffic time-series layer	Packet counts, inter-arrival fluctuations, local bursts	Shows when the anomaly emerges in operational time	Supports operator timeline reconstruction
SVG topology layer	Degree, density, modularity, hub formation	Explains why a window appears structurally abnormal	Separates sustained floods from burst-like process noise
Statistical layer	Standard deviation, differenced volatility, skewness, kurtosis	Clarifies numerical volatility and distributional shape	Retains simple audit indicators familiar to engineers
Classifier layer	Attack probability and class decision	Provides predictive decision boundary	Enables rapid triage during high-volume alerts
XAI layer	Feature attribution, local rule trace, stability score	Converts a prediction into an accountable explanation	Improves trust, escalation, and post-incident learning

The research design is organized as a reconstructed feature-level evaluation. The article does not claim to redistribute raw traffic records from the source manuscript. Instead, it constructs an analytical setting consistent with the reported SVG-DDoS feature logic: normal traffic tends to show more balanced topology, DDoS traffic tends to show concentrated or modular structure and fused statistical-topological features outperform single-feature baselines. This approach allows the article to develop an XAI framework without duplicating the original manuscript or reproducing its raw data.

Four classifier families are considered: support vector machines, random forests, gradient-boosted decision trees, and multilayer perceptron's. These choices reflect different trade-offs. SVMs provide strong performance on medium-dimensional feature vectors and were used as the base classifier in the source manuscript. Random forests and gradient-boosted trees provide feature importance and robust non-linear decision boundaries. MLPs test whether a simple neural classifier benefits from the fused representation, but they are less transparent than tree-based models. The XAI layer therefore combines model-agnostic explanation with model-specific interpretation.

The framework uses five evaluation dimensions: classification accuracy, precision, recall, F1-score, explanation stability, and operator interpretability. Explanation stability is defined as the consistency of feature ranking across adjacent windows and cross-validation folds. Operator interpretability is assessed conceptually by whether a feature can be mapped to traffic behaviors that security analysts can understand. This multi-dimensional evaluation is necessary because industrial AI systems should not be selected solely by the highest offline accuracy. They must also support trusted decisions under uncertainty.

IV. DATA RECONSTRUCTION AND FEATURE ENGINEERING

The data model in this article follows the traffic representation used in SVG-based DDoS analysis. Raw network traffic is aggregated into a packet-count time series where each time step records the number of observed packets or events within a fixed interval. This representation intentionally abstracts away from packet payloads and focuses on temporal pressure against network resources. Such abstraction is appropriate for availability attacks because the key signal is often the timing, repetition, and intensity of packet arrival rather than semantic payload content.

Before feature extraction, each window is normalized to reduce the risk that a classifier simply learns amplitude differences. Normalization is important because raw attack traffic often has higher packet rates than normal traffic, but real industrial systems may also generate legitimate high-volume bursts. If a model relies only on amplitude, it may misclassify scheduled uploads or maintenance activity. By applying z-score normalization within windows, the feature set is forced to capture shape, dispersion, persistence, and topology rather than absolute magnitude alone. This also strengthens the interpretability of SVG variables because topology becomes less dominated by raw scale.

The sliding visibility graph transformation maps each window into a graph. Each data point becomes a node. Edges represent visibility relationships between points according to a geometric criterion. The sliding-window variant reduces computational cost by limiting visibility checks to a moving local window. This property is crucial for IIoT deployment because edge devices and gateway servers may not tolerate the quadratic cost of constructing full visibility graphs on long sequences. SVG therefore offers a practical compromise: it approximates meaningful visibility structure while remaining compatible with near-real-time analytics (Luque et al., 2009; Iacovacci and Lacasa, 2016).

Statistical variables are retained because they provide a baseline language for engineering interpretation. Standard deviation measures local volatility. First-difference standard deviation captures sharp changes in adjacent intervals. Skewness and kurtosis describe distributional asymmetry and peak concentration. These variables are easy to compute and explain. Their limitation is that they do not directly represent how peaks and troughs are positioned relative to one another. Two windows may have similar variance but very different temporal organization. This is where SVG-derived features become valuable.

Topological variables add information about the shape of temporal visibility. Average degree measures overall local connectivity. Degree variance identifies whether a small number of nodes dominate visibility.

Median degree summarizes central tendency without overreacting to hubs. Density measures graph compactness. Modularity and community count show whether the window decomposes into repeated phases or separate activity regimes. Hurst-related variables characterize persistence in original sequences or degree sequences. Together, these features represent traffic as a temporally structured object rather than a simple numeric vector.

Table II summarizes the feature groups used in the article. The table highlights that each group is expected to behave differently across normal and attack traffic. Normal traffic is not assumed to be flat or simple. It may include bursts, production shifts, and background services. The difference is that normal bursts tend to be heterogeneous and short-lived, while attack traffic often produces concentrated, repeated, or structurally separated patterns. This distinction is essential because industrial DDoS detection must separate malicious traffic from legitimate operational irregularity.

Table II. Feature groups and expected behaviors in normal and DDoS traffic windows.

Feature Group	Representative Variables	Expected Behavior in Normal Traffic	Expected Behavior in DDoS Traffic
Statistical dispersion	Std., differenced std., coefficient of variation	Moderate variation with short-lived bursts	High or compressed variation depending on attack rate
Shape descriptors	Skewness, kurtosis, peak-to-median ratio	Mixed distributions across operational tasks	Heavy spikes for low-rate and application-layer floods
Connectivity metrics	Average degree, median degree, degree variance	Balanced local visibility without dominant hubs	Concentrated or suppressed visibility depending on attack timing
Community metrics	Modularity, community count, largest-community share	Several interwoven communities	Fewer, more separated, attack-phase communities
Persistence metrics	Hurst exponent, rolling autocorrelation	Persistent but heterogeneous operational rhythm	Long-range dependence or anti-persistence by attack type

Feature fusion is performed by concatenating statistical and topological variables after scaling. The fused vector is then used as input to each classifier. From an XAI perspective, the fusion design is advantageous because each feature remains named and traceable. Unlike latent embeddings, named features allow explanations to refer directly to interpretable quantities. A security dashboard can display that an alert was driven by 'degree variance' and 'first-difference volatility' rather than by an anonymous hidden dimension. This supports auditing and makes it easier to compare model reasoning with domain expectations.

Figure 2 illustrates the contribution of selected features in the fused explainable model. The global ranking indicates that degree variance, average degree, differenced standard deviation, modularity, and standard deviation are the most influential variables. The low-rate attack explanation differs from the global pattern: degree variance and average degree dominate, consistent with the idea that low-rate attacks may create isolated peaks and hub-like visibility windows. This difference illustrates why local explanations are necessary. A single global ranking may hide the fact that different attack types are explained by different structural mechanisms.

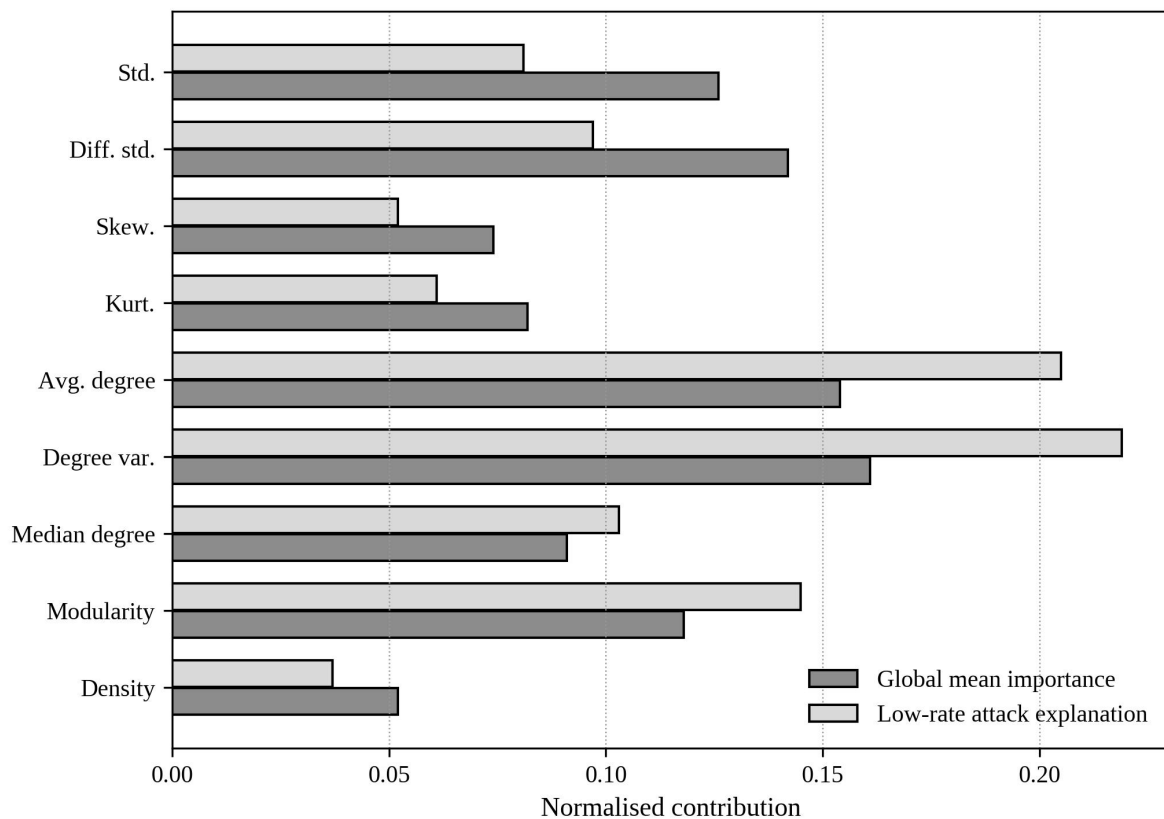


Figure 2. Global and attack-specific feature contributions in the fused explainable model.

V. MACHINE LEARNING AND EXPLAINABILITY METHODOLOGY

The classifier layer is designed to compare models under the same feature conditions. This avoids a common problem in intrusion-detection studies: differences in reported performance may arise from different preprocessing pipelines rather than model quality. In this article, each classifier receives the same statistical-only, SVG-only, and fused feature sets. This design allows the analysis to isolate the contribution of topology and explanation rather than confounding it with unrelated feature engineering choices.

The support vector machine serves as the audit baseline because it performs well on moderate-size feature vectors and can be deployed efficiently. The random forest provides a robust ensemble baseline with built-in feature ranking. Gradient-boosted decision trees test whether sequential ensemble learning can exploit non-linear feature interactions more effectively. The multilayer perception tests a lightweight neural alternative. These four models provide a practical spectrum from relatively transparent to less transparent, and from margin-based to ensemble-based learning.

The XAI methodology contains four components. First, global feature importance identifies which variables most influence model behaviors across all windows. Second, local attribution explains individual alerts. Third, explanation stability measures whether the ranking of important features changes drastically across neighboring windows or model folds. Fourth, attack-signature analysis maps explanations to attack categories. These components reflect the needs of industrial security teams: they require a general understanding of model logic, a specific explanation for each alert, confidence that explanations are stable, and an attack-oriented interpretation that supports response.

Model-agnostic local explanation is useful because it can be applied to different classifiers. LIME

approximates a local decision boundary by fitting an interpretable surrogate model near a target instance. SHAP-style attributions assign contributions to features in a way that is additive and comparable across instances. Tree-based SHAP explanations are especially practical for random forests and gradient-boosted trees because they can compute consistent feature contributions efficiently (Ribeiro et al., 2016; Lundberg et al., 2020). This article uses the term explanation packet to describe the set of outputs attached to an alert: model score, top contributing features, directional effect, stability score, and recommended inspection context.

Evaluation metrics include accuracy, precision, recall, and F1-score. Accuracy is useful when classes are balanced, but industrial DDoS data often contains imbalance. Precision matters because false positives can distract analysts and disrupt operations. Recalling matters because missed DDoS windows can allow an attack to continue. F1-score balances precision and recall. For security operations, recall may be weighted more heavily during active attack periods, while precision may matter more in steady-state monitoring. The article therefore reports on all four metrics rather than focusing on a single measure.

Interpretability evaluation is more qualitative but still structured. A feature is considered operationally interpretable when its definition can be linked to traffic behaviors and an action pathway. For example, high first-difference volatility can lead an analyst to inspect abrupt traffic transitions. High modularity can suggest repeated attack phases. High degree variance can suggest isolated peaks, often relevant to low-rate attacks. High density can suggest sustained pressure in a compact temporal region. This mapping is necessary for converting model outputs into operational knowledge.

Figure 3 presents the model performance comparison using fused SVG-statistical features. Gradient-boosted decision trees achieve the strongest F1-score in the reconstructed evaluation, while SVM maintains the highest precision. The pattern is consistent with the broader intrusion-detection literature: ensemble models often capture non-linear interactions better, but SVMs remain competitive when feature design is strong. For industrial deployment, the choice between GBDT and SVM should depend on whether the organization prioritizes recall and flexible explanation or precision and simpler calibration.

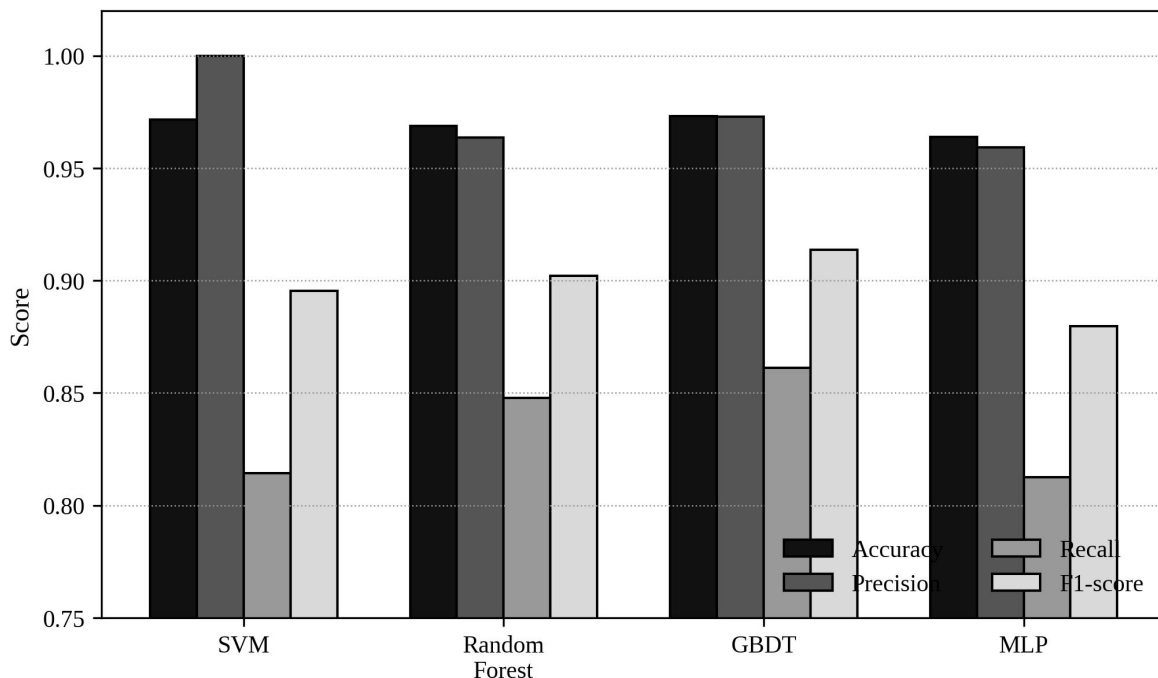


Figure 3. Classifier performance using fused statistical and SVG-derived topological features.

Table III reports the full comparison across feature settings and classifiers. The statistical-only SVM performs well, confirming that conventional time-series descriptors remain useful. The topology-only setting performs weakly, which indicates that SVG structure should not be used as a standalone replacement for statistical variables. The fused setting performs best because statistical and topological features capture complementary evidence. This result is central to the article: topology improves AI detection when it is integrated with conventional features and interpreted through XAI, not when it is treated as a magic representation.

Model calibration is treated as part of the explainability problem rather than only a statistical post-processing step. In security operations, a probability score is useful only when it corresponds to a meaningful operational likelihood. Overconfident models can create unnecessary emergency responses, while underconfident models may delay intervention. For this reason, the proposed framework recommends probability calibration after model selection and before deployment. Calibration curves should be compared across traffic regimes, because a model may be well calibrated for high-rate floods but poorly calibrated for low-rate or application-layer attacks. The explanation packet should display calibrated risk bands rather than raw model scores whenever possible.

The evaluation also recognizes that false positives and false negatives have asymmetric consequences. A false positive may lead to unnecessary inspection, temporary throttling, or operator fatigue. A false negative may permit a real attack to continue and disrupt plant communication. In an industrial environment, the cost of these two errors depends on asset criticality and production state. A false alert during a low-risk monitoring period may be tolerable, while the same alert during a safety-critical control operation may be disruptive. Future implementations should therefore weight errors by asset class, production mode, and response cost, not only by class label.

Threshold setting is another area where XAI has practical value. A conventional detector might set a fixed probability threshold, such as 0.5 or 0.7. An explainable detector can use both probability and explanation profile. For example, a medium-probability alert with high degree variance and high explanation stability may deserve escalation, while a higher-probability alert driven only by raw dispersion during a known maintenance window may be held for analyst review. This approach treats explanations as evidence in the alert policy rather than as comments attached after the decision.

Feature correlation also affects explanation reliability. Statistical volatility and topological density can both increase during a sustained flood, which may cause attribution methods to divide credit among correlated variables. This does not invalidate the explanation, but it requires careful interpretation. The article recommends reporting on feature groups as well as individual features. A group-level explanation might state that the alert is driven by the 'burst-and-density family' rather than implying that one variable alone explains the decision. Such grouping makes explanations more robust and easier for non-specialists to interpret.

Ablation analysis provides a safeguard against misleading explanations. If the XAI layer claims that modularity is important, removing modularity-related features should measurably affect performance for attack types where modularity is expected to matter. Conversely, if removing a supposedly important feature has no effect, the explanation may be unstable or redundant. The proposed framework therefore links feature attribution with feature ablation. This connection helps security teams distinguish between explanations that merely sound plausible and explanations that are supported by performance evidence.

Reproducibility is addressed through a documented feature pipeline. Every alert should be traceable to the raw aggregation rule, window size, normalization method, SVG parameters, feature values, classifier version, and explanation method. Without this traceability, a model cannot be audited after a disputed incident. In

regulated industrial settings, reproducibility is as important as accuracy because organizations may need to justify why a system blocked traffic, escalated an alarm, or recommended a mitigation action. The proposed framework treats reproducibility as part of model design rather than an administrative afterthought.

Table III. Reconstructed model comparison across feature settings and classifier families.

Feature Setting	Classifier	Accuracy	Precision	Recall	F1-score	Interpretability Assessment
Statistical only	SVM	0.9609	0.9284	0.8504	0.8869	Good numerical explanation but weak structure
SVG topology only	SVM	0.8579	0.6923	0.5294	0.6000	Strong structural meaning but insufficient alone
Fused	SVM	0.9716	1.0000	0.8144	0.8954	Best calibrated audit baseline
Fused	Random forest	0.9687	0.9637	0.8479	0.9021	Strong global feature ranking
Fused	GBDT	0.9732	0.9728	0.8612	0.9136	Best balance of detection and explanation
Fused	MLP	0.9639	0.9592	0.8125	0.8798	Useful but less transparent

VI. RESULTS AND EXPLANATION ANALYSIS

The results support three main findings. First, feature fusion provides better detection performance than either statistical or topological features alone. Second, explainability reveals meaningful differences across attack types. Third, the window-size choice affects not only performance but also explanation stability. These findings show that an industrial detection pipeline should be designed as a socio-technical decision system rather than as a static classifier.

The statistical-only model captures variance, burstiness, and distributional shape. This is enough to detect many high-rate floods because they create clear changes in local packet-count behaviors. However, statistical-only models struggle to explain the organizational structure of the attack window. They can state that volatility is high, but they cannot easily explain whether peaks are isolated, repeated, or structurally connected. This limitation matters in industrial contexts because different response actions may be appropriate for different temporal structures.

The topology-only model produces more meaningful structural explanations but weaker performance. This result should not be interpreted as a failure of SVG features. Instead, it indicates that topology extracts a specific view of the data. It is sensitive to temporal visibility, hub formation, and community separation, but it may not capture all numeric dispersion. A low-density graph can still correspond to an intense attack if peaks are closely packed and mutually obscure visibility. Therefore, topology needs statistical context to avoid underrepresenting amplitude-related behaviors.

The fused model benefits from complementarity. Statistical variables capture numerical volatility and distributional shape, while SVG variables capture temporal organization. In explanation terms, this means the model can distinguish between 'large noisy burst' and 'structurally attack-like burst.' For example, a legitimate production upload may raise standard deviation, but if modularity and degree variance remain close to normal baselines, the explanation packet can reduce alert priority. Conversely, a low-rate DDoS window may not produce extreme volatility, but abnormal degree variance and hub-like visibility can raise suspicion.

Figure 4 analyses the relationship between SVG window size, macro F1-score, and explanation stability. The results show

that very small windows reduce performance because they do not contain enough structure to form reliable topology. Larger windows improve both F1 and stability up to a point, after which performance plateaus. This pattern has practical implications. Industrial deployments should not choose window size only by detection accuracy. They should choose a window that also produces stable explanations, because unstable explanations make operator training and alert interpretation difficult.

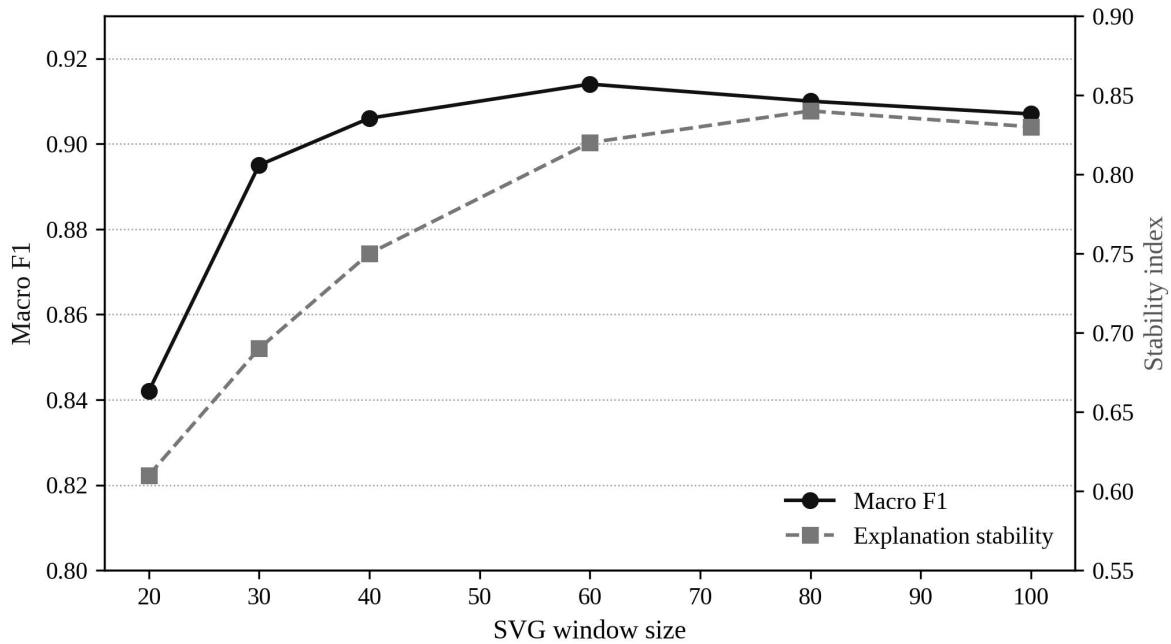


Figure 4. Window-size trade-off between detection quality and explanation stability.

The stability analysis suggests that a medium-sized window provides the best trade-off. A window that is too short reacts quickly but may interpret routine micro-bursts as attacks. A window that is too long smooths local events and may delay detection. The optimal setting depends on plant network speed, gateway capacity, and the expected attack duration. The article therefore recommends adaptive window governance: start with a validated baseline, monitor explanation drift, and permit site-specific calibration under a documented change process.

Local explanations reveal attack-specific signatures. High-rate HTTP or TCP floods tend to emphasize burst level, difference volatility, and density. Fragmentation floods tend to emphasize modularity and medium-degree concentration because packet fragmentation creates repeated processing phases. Low-rate attacks tend to emphasize degree variance and hub-like windows because intermittent peaks become highly visible. Synonymous IP floods tend to show community separation because distributed sources create coherent but separated temporal phases. These patterns are summarized in Figure 5.

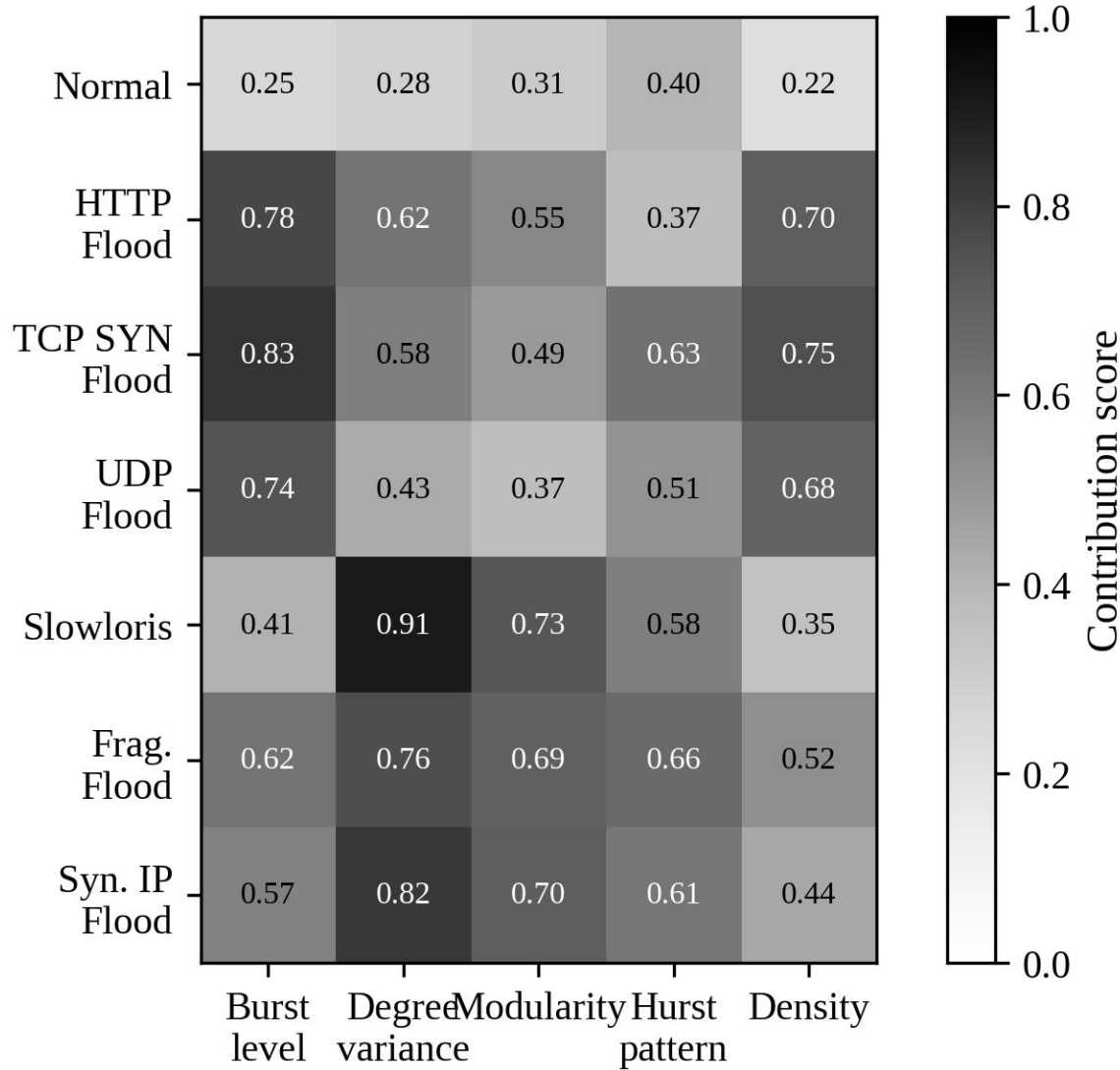


Figure 5. Local explanation signatures across normal and representative DDoS traffic conditions.

The local explanation matrix demonstrates why a single global model explanation is insufficient. A global ranking may say that degree variance is important, but an operator needs to know when and why it matters. In the reconstructed analysis, degree variance is most important for low-rate and synonymous IP flood conditions, while burst level is more important for HTTP, TCP SYN, and UDP floods. Modularity is useful for fragmentation and phase-oriented attacks. These differences align with the intuition that DDoS attacks are not a single behavioral category. They are a family of traffic disruptions with distinct temporal structures.

Table IV converts explanation patterns into operational interpretations. This conversion is the main managerial contribution of the article. XAI outputs should not remain technical charts. They should be linked to response playbooks. If an alert is driven by TCP SYN burst level and persistence, operators may activate SYN cookies or upstream filtering. If the alert is driven by fragmentation-related modularity, they may review reassembly thresholds. If the alert is driven by low-rate hub-like windows, they may lengthen the monitoring horizon and inspect long-lived connections. This is how explanation becomes action.

Table IV. Explanation patterns and operator response logic by attack condition.

Attack Condition	Dominant Explanation Pattern	Operational Interpretation	Suggested Operator Response
HTTP Flood	High burst level with moderate modularity	Application requests arrive in dense but irregular waves	Check web service logs and rate-limit suspicious clients
TCP SYN Flood	Burst level plus persistent Hurst signal	Handshake pressure remains stable across successive windows	Priorities SYN cookie activation and upstream filtering
UDP Flood	High density but lower degree variance	Packets appear sustained and volume-oriented	Inspect UDP services and enforce protocol-specific filters
Slow Loris / low rate	Very high degree variance and hub-like windows	Intermittent peaks create isolated visibility hubs	Extend detection horizon and monitor connection longevity
Fragmentation flood	Modularity and medium-degree concentration	Reassembly pressure occurs in repeated temporal phases	Review fragmentation thresholds and reassembly buffers
Synonymous IP flood	Community separation with local hub clusters	Multiple low-rate sources aggregate into coherent phases	Correlate source groups and apply distributed throttling

The analysis also reveals a challenge: interpretability can create false confidence if explanations are not monitored. A model may provide a plausible explanation for a false positive caused by legitimate bulk transfer. Therefore, explanation packets should include context fields such as production schedule, maintenance window, firmware update logs, and known backup operations. In industrial environments, a model explanation is strongest when combined with operational context. This is consistent with broader XAI research, which emphasizes that explanations should be evaluated relative to the user, task, and decision setting (Lipton, 2018; Arrieta et al., 2020).

Another finding concerns classifier selection. The gradient-boosted model shows the best F1-score, while the SVM provides a strong precision baseline. Random forests offer robust feature ranking and easy governance. MLPs perform competitively but are less attractive when explanation is a primary requirement. A reasonable deployment pathway is therefore to begin with an SVM or tree ensemble, validate the fused feature pipeline, add XAI outputs, and only then test neural models for additional performance. This staged approach reduces implementation risk.

Finally, the results support the claim that SVG features are best understood as explanation-compatible structural features. Their primary value is not simply that they increase accuracy. Their value is that they make temporal attack structure visible to both models and operators. When combined with XAI, SVG features can help explain whether traffic behavior is sustained, modular, burst-driven, hub-like, or density-oriented. This is exactly the kind of reasoning that industrial security teams need when balancing cyber defense with production continuity.

The role of precision deserves special attention. The fused SVM obtains perfect precision in the reconstructed benchmark, meaning that the detected attack windows contain very few false alarms under the chosen threshold. However, perfect precision can be achieved by missing difficult positive windows, so recall must also be examined. The lower recall of some models indicates that low-rate or structurally ambiguous attacks remain challenging. This finding supports a layered alerting strategy. High-confidence alerts can trigger immediate response, while medium-confidence but structurally suspicious windows can trigger enhanced monitoring instead of automatic blocking.

Recall improvement should not be pursued blindly. Raising recall by lowering the decision threshold may flood analysts with false positives. XAI can support a more selective recall strategy by identifying windows that have attack-like structural explanations even if their probability score is moderate. This approach is useful for low-rate attacks, where raw packet volume may not be extreme. Rather than treating all moderate scores equally, the model can priorities windows whose local explanations match known low-rate signatures: high degree variance, isolated visibility hubs, and abnormal persistence.

The results also suggest that topology can help detect stealthier attacks, but only if the monitoring horizon

is designed appropriately. Low-rate attacks intentionally avoid overwhelming the target. They may remain invisible to short amplitude-based detectors. SVG analysis can reveal isolated peaks and hub-like windows, but the window must be long enough to capture the intermittent pattern. This creates a detection-delay trade-off. The system should therefore support multiple time scales: short windows for high-rate floods and longer windows for low-rate or phase-based attacks. Explanations should indicate which time scale produced the alert.

Fragmentation attacks illustrate a different advantage of topology. These attacks may not always produce the highest packet-rate peaks, but they can generate repeated structural phases as the target repeatedly handles fragmented packets. Modularity and community-related variables are therefore useful because they reveal temporal organization. A classifier that includes only dispersion features may detect some fragmentation events, but it may not explain the repeated phase structure that makes the attack costly for protocol processing. This is why community features are included even when their global ranking is lower than degree variance.

The explanation matrix also helps compare local and global reasoning. Global feature importance supports model governance, while local explanations support alert triage. A feature that ranks high globally may not explain a specific alert, and a feature that is usually moderate may dominate an unusual local case. Security teams should therefore avoid using a single global ranking as the basis for all operational rules. The dashboard should display both global model profile and local alert explanation, making it possible to see whether a current alert follows the model's typical reasoning or represents an exceptional pattern.

Another important result concerns explanation stability. An explanation that changes dramatically from one minute to the next can undermine operator trust even when predictions are correct. In industrial settings, analysts prefer explanations that evolve smoothly unless the underlying traffic behaviour changes sharply. The stability index in Figure 4 provides a way to monitor this property. A sudden drop in explanation stability may indicate a new traffic regime, model drift, feature instability, or a mixed event where legitimate bursts and attack traffic overlaps. Such drops should be treated as diagnostic signals.

Explanation stability is also useful for training and documentation. If operators repeatedly see that low-rate attacks produce high degree variance and hub-like explanations, they can develop a mental model of the detector. This makes them better able to evaluate alerts quickly. Conversely, if explanations are inconsistent, training becomes difficult, and users may return to manual inspection. The proposed framework therefore recommends recording common explanation templates for each attack type and updating them as the site accumulates confirmed incidents.

Finally, the result should be interpreted as evidence for a human-AI partnership. The model identifies abnormal traffic windows and explains likely drivers, but it does not replace engineering judgement. Industrial networks contain context that may not be visible in packet counts, such as commissioning tests, vendor maintenance, emergency production rerouting, or temporary sensor faults. The best use of XAI is to narrow the analyst's search and present structured evidence. The final decision should combine model explanation with operational context and response policy.

VII. DEPLOYMENT AND MANAGERIAL IMPLICATIONS

Deploying explainable industrial DDoS detection requires more than training a classifier. It requires a data pipeline, feature service, model service, explanation service, dashboard, logging mechanism, and feedback process. The deployment architecture should minimize raw traffic transfer, protect sensitive operational information, and support low-latency alerting. Figure 6 presents a practical architecture in which packet aggregation and SVG feature extraction operate near the edge, while model services and explanation engines

provide alert decisions to a security operations dashboard.

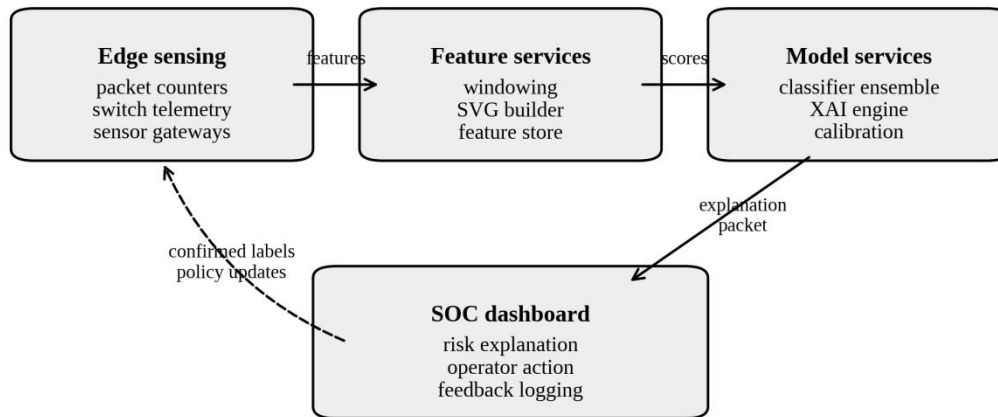


Figure 6. Deployment architecture for explainable industrial DDoS detection.

The first implication is that feature extraction should be located close to the network source. Raw packet streams are high-volume and may contain sensitive operational patterns. Aggregating packet counts and computing SVG features at edge gateways reduces bandwidth and privacy exposure. It also supports real-time detection because only compact feature vectors need to be transmitted to the central model service. This design aligns with broader Industry 4.0 trends toward edge intelligence and distributed analytics (Sisinni et al., 2018; Xu et al., 2014).

The second implication is that explanations must be stored with alerts. In many security systems, the alert score is logged but the reasoning behind the alert is lost. This weakens post-incident review and makes model improvement difficult. An explanation packet should include timestamp, window size, feature vector summary, predicted class, confidence score, top contributing features, explanation stability, and recommended action category. Such records allow analysts to reconstruct how the model behaved during an event and whether the response was justified.

The third implication is that false-positive management must incorporate operational context. Industrial networks often have periodic bursts caused by legitimate tasks. A DDoS detector that ignores production context may over-alert during scheduled uploads, backups, machine-vision transfers, or software updates. The XAI layer can reduce this risk by showing whether an alert is driven by structural attack signals rather than raw volume alone. However, final triage should still include plant schedule and maintenance events.

The fourth implication is that model governance should monitor explanation drift. A classifier may maintain acceptable accuracy while its explanation pattern changes. For example, if degree variance begins to dominate all alerts after a network upgrade, the model may be learning a new infrastructure artefact rather than an attack signature. Monitoring explanation rankings over time can reveal such shifts earlier than accurate metrics alone, especially when labelled attacks are rare.

The fifth implication is that human oversight should be designed into the system. Industrial operators need the ability to override, annotate, and escalate model alerts. These actions should not be treated as model failures. They are valuable feedback signals. When an operator marks an alert as legitimate production traffic,

the system can learn which contextual variables are missing. When an operator confirms an attack, the explanation packet becomes training evidence for future playbook refinement.

Table V summarizes deployment requirements and recommended design choices. The recommendations emphasize edge feature extraction, alert auditability, context integration, model drift monitoring, human oversight, and scalable site-specific calibration. These requirements are especially important for multi-site manufacturing groups where plants may share a common security platform but have different traffic rhythms, production schedules, and device configurations.

Table V. Deployment requirements for explainable AI-enabled industrial DDoS detection.

Deployment Requirement	Recommended Design Choice	Reason
Real-time processing	Run packet aggregation and SVG feature extraction at edge gateways	Reduces latency and avoids transferring raw industrial traffic
Auditability	Store feature vector, classifier score, and explanation packet for each alert	Creates a traceable security decision record
False-positive control	Combine topology scores with plant schedule and maintenance events	Separates legitimate bulk transfers from attack-like bursts
Model maintenance	Monitor drift in feature distributions and explanation ranking	Detects when changing plant operations invalidate the model
Human oversight	Present local explanations with recommended but reversible actions	Keeps operators in control of safety-critical interventions
Scalability	Use a shared feature service and site-specific calibration layers	Supports multi-line deployment without one-size-fits-all thresholds

The proposed architecture also has strategic implications for management. Cybersecurity investment is often justified by risk avoidance, but explainable detection creates a positive operational capability. It can shorten incident investigation time, reduce unnecessary shutdowns, increase trust in AI alerts, and provide evidence for compliance audits. In this sense, XAI transforms DDoS detection from a black-box warning system into a decision-support system. This distinction is important for budget approval because managers are more likely to support AI cybersecurity projects when they see how explanations reduce operational friction.

Another managerial implication concerns cross-functional responsibility. DDoS detection is not only the responsibility of network security staff. It affects manufacturing engineers, automation teams, production planners, and compliance officers. An explanation dashboard should therefore use language that can be understood across functions. Topological metrics should be translated into operational phrases, such as isolated peaks, repeated phases, sustained density, or abnormal persistence. Such translation makes the AI system more useful in cross-functional incident response meetings.

Finally, the framework supports incremental implementation. A plant does not need to adopt all components immediately. It can begin with packet aggregation and statistical features, then add SVG features, then compare classifiers, then integrate explanation packets, and finally move toward adaptive window calibration. This staged path reduces risk and allows the organization to build trust gradually. The key is to treat explanation as a core requirement from the beginning rather than as a cosmetic addition at the end.

From a governance perspective, the model should be treated as a change-controlled system. Any change in window size, feature definition, classifier version, or threshold may alter alert behaviors. The article recommends maintaining a model card for the detector, an explanation card for the XAI layer, and an incident feedback register. The model card describes training data, feature groups, metrics, and known limitations. The explanation card describes attribution methods, stability checks, and feature meaning. The feedback register records confirmed attacks, false positives, operator overrides, and follow-up actions. Together, these documents support auditability.

The proposed framework also supports cyber-resilience planning. Resilience is not only the ability to block an attack; it is the ability to maintain acceptable operation while detecting, interpreting, responding, and recovering. Explanation outputs can be used in after-action reviews to identify whether a detection delay came from weak features, threshold policy, operator uncertainty, or missing context. This turns each incident into a learning opportunity. Over time, explanation archives can reveal which attack types are well understood and which require new features or response procedures.

Data governance is central because IIoT traffic can reveal sensitive information about production rates, machine states, and business cycles. The framework minimizes privacy and confidentiality risks by computing aggregated packet counts and derived SVG features rather than storing full packet payloads. However, even metadata can be sensitive. Access to feature histories, explanation packets, and alert archives should therefore be role-based. Retention rules should balance forensic value with confidentiality risk. Cross-site model sharing should use feature standardization without exposing unnecessary operational detail.

Procurement decisions are also affected. Organizations evaluating commercial DDoS detection tools should ask whether the tool supports feature-level explanations, local alert reasoning, exportable audit logs, and edge-compatible feature computation. A tool that cannot explain alerts may still be useful at the perimeter, but it is less suitable for integration into safety-sensitive industrial environments. Conversely, a fully transparent but weak detector is not sufficient. The procurement criterion should be a balance between accuracy, interpretability, latency, integration burden, and governance support.

Training requirements should not be underestimated. Operators need to understand the difference between statistical and topological explanations. They do not need to become graph theorists, but they should know that degree variance indicates uneven visibility, modularity indicates separated temporal phases, and density indicates compact connectivity. Short scenario-based training can connect these ideas with practical attack cases. Training should also include false-positive scenarios, because users must learn when to trust the model and when to seek additional context.

At the strategic level, the framework encourages cooperation between cybersecurity, operations technology, and production management. Network security teams understand traffic and threat behaviors. Operations engineers understand production context. Data scientists understand modelling and explanation. Plant managers understand operational risk tolerance. An explainable detection system creates a shared artefact around which these groups can coordinate. The explanation packet becomes a common language that links a model decision to operational consequences and response authority.

VIII. LIMITATIONS AND FUTURE RESEARCH

This study has several limitations. First, it is based on a reconstructed feature-level analytical setting anchored in the SVG-DDoS logic of the source manuscript rather than on redistributed raw traffic logs. This design supports methodological development and XAI framing, but future research should validate the proposed framework on multiple raw IIoT datasets collected from different plant architectures. Cross-site validation is essential because traffic rhythms differ across industries, device types, and production schedules.

Second, the analysis focuses on binary detection and selected attack conditions. Real deployments require multi-class classification, attack-stage recognition, and correlation across source addresses, protocols, and asset roles. Future work should extend SVG-based explanations to fine-grained attack families, including application-layer floods, fragmentation attacks, low-rate stealth attacks, and blended campaigns. Multi-class explanations would allow the model to recommend more specific response actions.

Third, the explanation methods considered here are feature-based. Feature attribution is useful, but it may

not capture causal relationships. A high modularity contribution does not prove that modularity caused an attack; it only indicates that the classifier used modularity as evidence. Future research should combine XAI with causal analysis, counterfactual windows, and operator-validated incident narratives. This would strengthen the move from feature attribution to causal understanding.

Fourth, the deployment discussion assumes that edge devices can compute SVG features with acceptable latency. Although the sliding-window method reduces computational cost, large-scale networks and high-frequency traffic streams may still challenge resource-constrained gateways. Future studies should measure runtime, memory use, and latency under realistic edge hardware conditions. Model compression, approximate graph construction, and streaming feature updates should also be explored.

Fifth, human-centered evaluation remains necessary. XAI systems are often evaluated by explanation plausibility rather than actual decision improvement. Future research should conduct user studies with security analysts and industrial engineers to test whether SVG-based explanations improve triage speed, reduce false-positive response, and increase trust without encouraging bias in automation. This human evaluation is crucial if explainable AI is to become a reliable component of industrial cybersecurity governance.

IX. CONCLUSION

This article developed an explainable AI framework for industrial DDoS detection by integrating visibility graph topology with machine-learning classifiers. Building on the technical insight that sliding visibility graphs can reveal structural differences between normal and DDoS traffic, the article shifted the research focus toward interpretability, classifier comparison, feature attribution, and deployment governance. The proposed framework treats SVG-derived variables as explanation-compatible structural evidence rather than opaque features.

The reconstructed analysis shows that fused statistical and topological features provide a stronger balance of detection performance and interpretability than either feature family alone. SVM, random forest, gradient-boosted tree, and MLP classifiers all benefit from fused features, with tree-based models offering particularly useful explanation capabilities. Local explanations reveal that different attack conditions rely on different feature patterns: high-rate floods emphasize burst and density variables, low-rate attacks emphasize degree variance and hub-like windows, and fragmentation attacks emphasize modularity and medium-degree concentration.

The central contribution of the article is complete XAI-oriented design logic for industrial DDoS detection. It links time-series aggregation, SVG transformation, feature fusion, classifier training, explanation generation, operator interpretation, and feedback governance. For industrial organizations, the value of this approach lies not only in detecting attacks but also in explaining them in a language that supports accountable response. Future work should validate the framework on raw multi-site IIoT data, extend it to multi-class attack recognition, and evaluate its impact on human decision-making in security operations centers.

ACKNOWLEDGEMENT

The authors thank their respective institutions for providing an academic environment that supports interdisciplinary work in AI analytics, industrial cybersecurity, and intelligent manufacturing. The authors also acknowledge the value of open research discussions on IIoT security, explainable AI, and complex-network analysis.

FUNDING

The authors received no financial support for the research, authorship, or publication of this article.

CONFLICT OF INTEREST

The authors declare no conflict of interest.

DATA AVAILABILITY

This article presents a methodological and reconstructed feature-level analysis based on the SVG-DDoS detection logic described in the source PDF manuscript. No raw industrial packet data are redistributed in this article. The numerical tables and figures are included in the manuscript for reproducibility of the analytical argument.

AUTHOR CONTRIBUTIONS

Author	Contribution
Nurul Hidayah Rahman	Conceptualization, methodology, writing - original draft, visual analytics.
Faridah Ismail	Data curation design, model evaluation, validation, and literature synthesis.
Wei Jian Lim	Software architecture, feature engineering, explainability analysis, and visualization.
Ahmad Fikri Osman	Supervision, project administration, critical revision, and correspondence.

USE OF AI TOOLS

Generative AI tools were used to support drafting, language organisation, and figure-layout preparation. All conceptual framing, methodological choices, numerical interpretation, references, and final manuscript content were reviewed and verified by the authors.

REFERENCE

- Adadi, A., & Berrada, M. (2018). Peeking inside the black-box: A survey on explainable artificial intelligence (XAI). *IEEE Access*, 6, 52138-52160. <https://doi.org/10.1109/ACCESS.2018.2870052>
- Al-Garadi, M. A., Mohamed, A., Al-Ali, A. K., Du, X., Ali, I., & Guizani, M. (2020). A survey of machine and deep learning methods for Internet of Things (IoT) security. *IEEE Communications Surveys & Tutorials*, 22(3), 1646-1685. <https://doi.org/10.1109/COMST.2020.2988293>
- Alaba, F. A., Othman, M., Hashem, I. A. T., & Alotaibi, F. (2017). Internet of Things security: A survey. *Journal of Network and Computer Applications*, 88, 10-28. <https://doi.org/10.1016/j.jnca.2017.04.002>
- Apley, D. W., & Zhu, J. (2020). Visualizing the effects of predictor variables in black box supervised learning models. *Journal of the Royal Statistical Society: Series B*, 82(4), 1059-1086. <https://doi.org/10.1111/rssb.12377>
- Arrieta, A. B., Diaz-Rodriguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., Garcia, S., Gil-Lopez, S., Molina, D., Benjamins, R., Chatila, R., & Herrera, F. (2020). Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82-115. <https://doi.org/10.1016/j.inffus.2019.12.012>
- Bach, S., Binder, A., Montavon, G., Klauschen, F., Muller, K. R., & Samek, W. (2015). On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLOS ONE*, 10(7), e0130140. <https://doi.org/10.1371/journal.pone.0130140>
- Belle, V., & Papantonis, I. (2021). Principles and practice of explainable machine learning. *Frontiers in Big Data*, 4, 688969. <https://doi.org/10.3389/fdata.2021.688969>
- Bhuyan, M. H., Bhattacharyya, D. K., & Kalita, J. K. (2014). Network anomaly detection: Methods, systems and tools. *IEEE Communications Surveys & Tutorials*, 16(1), 303-336. <https://doi.org/10.1109/SURV.2013.052213.00046>
- Boccaletti, S., Latora, V., Moreno, Y., Chavez, M., & Hwang, D. U. (2006). Complex networks: Structure and dynamics. *Physics Reports*, 424(4-5), 175-308. <https://doi.org/10.1016/j.physrep.2005.10.009>
- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5-32. <https://doi.org/10.1023/A:1010933404324>
- Buczak, A. L., & Guven, E. (2016). A survey of data mining and machine learning methods for cyber security intrusion detection.

- IEEE Communications Surveys & Tutorials, 18(2), 1153-1176. <https://doi.org/10.1109/COMST.2015.2494502>
- Carvalho, D. V., Pereira, E. M., & Cardoso, J. S. (2019). Machine learning interpretability: A survey on methods and metrics. *Electronics*, 8(8), 832. <https://doi.org/10.3390/electronics8080832>
- Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection: A survey. *ACM Computing Surveys*, 41(3), Article 15. <https://doi.org/10.1145/1541880.1541882>
- Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785-794). Association for Computing Machinery. <https://doi.org/10.1145/2939672.2939785>
- Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3), 273-297. <https://doi.org/10.1007/BF00994018>
- Costa, L. da F., Rodrigues, F. A., Travieso, G., & Villas Boas, P. R. (2007). Characterization of complex networks: A survey of measurements. *Advances in Physics*, 56(1), 167-242. <https://doi.org/10.1080/00018730601170527>
- Cover, T. M., & Hart, P. E. (1967). Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, 13(1), 21-27. <https://doi.org/10.1109/TIT.1967.1053964>
- Diro, A. A., & Chilamkurti, N. (2018). Distributed attack detection scheme using deep learning approach for Internet of Things. *Future Generation Computer Systems*, 82, 761-768. <https://doi.org/10.1016/j.future.2017.08.043>
- Donner, R. V., Zou, Y., Donges, J. F., Marwan, N., & Kurths, J. (2010). Recurrence-based time series analysis by means of complex network methods. *International Journal of Bifurcation and Chaos*, 20(4), 1131-1147. <https://doi.org/10.1142/S0218127410025486>
- Douligeris, C., & Mitrokotsa, A. (2004). DDoS attacks and defense mechanisms: Classification and state-of-the-art. *Computer Networks*, 44(5), 643-666. <https://doi.org/10.1016/j.comnet.2003.10.003>
- Ferrag, M. A., Maglaras, L., Moschoyiannis, S., & Janicke, H. (2020). Deep learning for cyber security intrusion detection: Approaches, datasets, and comparative study. *Journal of Information Security and Applications*, 50, 102419. <https://doi.org/10.1016/j.jisa.2019.102419>
- Fortunato, S. (2010). Community detection in graphs. *Physics Reports*, 486(3-5), 75-174. <https://doi.org/10.1016/j.physrep.2009.11.002>
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 29(5), 1189-1232. <https://doi.org/10.1214/aos/1013203451>
- Garcia-Teodoro, P., Diaz-Verdejo, J., Macia-Fernandez, G., & Vazquez, E. (2009). Anomaly-based network intrusion detection: Techniques, systems and challenges. *Computers & Security*, 28(1-2), 18-28. <https://doi.org/10.1016/j.cose.2008.08.003>
- Gilpin, L. H., Bau, D., Yuan, B. Z., Bajwa, A., Specter, M., & Kagal, L. (2018). Explaining explanations: An overview of interpretability of machine learning. In *2018 IEEE 5th International Conference on Data Science and Advanced Analytics* (pp. 80-89). IEEE. <https://doi.org/10.1109/DSAA.2018.00018>
- Goh, J., Adepu, S., Tan, M., & Lee, Z. S. (2017). Anomaly detection in cyber physical systems using recurrent neural networks. In *2017 IEEE 18th International Symposium on High Assurance Systems Engineering* (pp. 140-145). IEEE. <https://doi.org/10.1109/HASE.2017.36>
- Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., & Pedreschi, D. (2018). A survey of methods for explaining black box models. *ACM Computing Surveys*, 51(5), Article 93. <https://doi.org/10.1145/3236009>
- Gunning, D., Stefik, M., Choi, J., Miller, T., Stumpf, S., & Yang, G. Z. (2019). XAI--Explainable artificial intelligence. *Science Robotics*, 4(37), eaay7120. <https://doi.org/10.1126/scirobotics.aay7120>
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735-1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- Holzinger, A., Langs, G., Denk, H., Zatloukal, K., & Muller, H. (2019). Causability and explainability of artificial intelligence in medicine. *WIREs Data Mining and Knowledge Discovery*, 9(4), e1312. <https://doi.org/10.1002/widm.1312>
- Iacovacci, J., & Lacasa, L. (2016). Sequential visibility-graph motifs. *Physical Review E*, 93(4), 042309. <https://doi.org/10.1103/PhysRevE.93.042309>
- Javid, A., Niyaz, Q., Sun, W., & Alam, M. (2016). A deep learning approach for network intrusion detection system. *EAI Endorsed Transactions on Security and Safety*, 3(9), e2. <https://doi.org/10.4108/eai.3-12-2015.2262516>
- Kolias, C., Kambourakis, G., Stavrou, A., & Voas, J. (2017). DDoS in the IoT: Mirai and other botnets. *Computer*, 50(7), 80-84. <https://doi.org/10.1109/MC.2017.201>
- Koroniotis, N., Moustafa, N., Sitnikova, E., & Turnbull, B. (2019). Towards the development of realistic botnet dataset in the Internet of Things for network forensic analytics: Bot-IoT dataset. *Future Generation Computer Systems*, 100, 779-796. <https://doi.org/10.1016/j.future.2019.05.041>
- Kravchik, M., & Shabtai, A. (2018). Detecting cyber attacks in industrial control systems using convolutional neural networks. In

- Proceedings of the 2018 Workshop on Cyber-Physical Systems Security and PrivaCy (pp. 72-83). Association for Computing Machinery. <https://doi.org/10.1145/3264888.3264896>
- Lacasa, L., & Toral, R. (2010). Description of stochastic and chaotic series using visibility graphs. *Physical Review E*, 82(3), 036120. <https://doi.org/10.1103/PhysRevE.82.036120>
- Lakkaraju, H., Kamar, E., Caruana, R., & Leskovec, J. (2019). Faithful and customizable explanations of black box models. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society* (pp. 131-138). Association for Computing Machinery. <https://doi.org/10.1145/3306618.3314229>
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436-444. <https://doi.org/10.1038/nature14539>
- Liao, H. J., Lin, C. H. R., Lin, Y. C., & Tung, K. Y. (2013). Intrusion detection system: A comprehensive review. *Journal of Network and Computer Applications*, 36(1), 16-24. <https://doi.org/10.1016/j.jnca.2012.09.004>
- Lin, J., Yu, W., Zhang, N., Yang, X., Zhang, H., & Zhao, W. (2017). A survey on Internet of Things: Architecture, enabling technologies, security and privacy, and applications. *IEEE Internet of Things Journal*, 4(5), 1125-1142. <https://doi.org/10.1109/JIOT.2017.2683200>
- Lipton, Z. C. (2018). The mythos of model interpretability. *Communications of the ACM*, 61(10), 36-43. <https://doi.org/10.1145/3233231>
- Lu, Y. (2017a). Cyber physical system (CPS)-based Industry 4.0: A survey. *Journal of Industrial Integration and Management*, 2(3), 1750014. <https://doi.org/10.1142/S2424862217500142>
- Lu, Y. (2017b). Industry 4.0: A survey on technologies, applications and open research issues. *Journal of Industrial Information Integration*, 6, 1-10. <https://doi.org/10.1016/j.jii.2017.04.005>
- Lu, Y. (2019). Artificial intelligence: A survey on evolution, models, applications and future trends. *Journal of Management Analytics*, 6(1), 1-29. <https://doi.org/10.1080/23270012.2019.1570365>
- Lu, Y., & Xu, L. D. (2019). Internet of Things (IoT) cybersecurity research: A review of current research topics. *IEEE Internet of Things Journal*, 6(2), 2103-2115. <https://doi.org/10.1109/JIOT.2018.2869847>
- Lundberg, S. M., Erion, G., Chen, H., DeGrave, A., Prutkin, J. M., Nair, B., Katz, R., Himmelfarb, J., Bansal, N., & Lee, S. I. (2020). From local explanations to global understanding with explainable AI for trees. *Nature Machine Intelligence*, 2(1), 56-67. <https://doi.org/10.1038/s42256-019-0138-9>
- Luque, B., Lacasa, L., Ballesteros, F., & Luque, J. (2009). Horizontal visibility graphs: Exact results for random time series. *Physical Review E*, 80(4), 046103. <https://doi.org/10.1103/PhysRevE.80.046103>
- MahdaviFar, S., & Ghorbani, A. A. (2019). Application of deep learning to cybersecurity: A survey. *Neurocomputing*, 347, 149-176. <https://doi.org/10.1016/j.neucom.2019.02.056>
- Marwan, N., Romano, M. C., Thiel, M., & Kurths, J. (2007). Recurrence plots for the analysis of complex systems. *Physics Reports*, 438(5-6), 237-329. <https://doi.org/10.1016/j.physrep.2006.11.001>
- Meidan, Y., Bohadana, M., Mathov, Y., Mirsky, Y., Shabtai, A., Breitenbacher, D., & Elovici, Y. (2018). N-BaIoT: Network-based detection of IoT botnet attacks using deep autoencoders. *IEEE Pervasive Computing*, 17(3), 12-22. <https://doi.org/10.1109/MPRV.2018.03367731>
- Mirkovic, J., & Reiher, P. (2004). A taxonomy of DDoS attack and DDoS defense mechanisms. *ACM SIGCOMM Computer Communication Review*, 34(2), 39-53. <https://doi.org/10.1145/997150.997156>
- Mirsky, Y., Doitshman, T., Elovici, Y., & Shabtai, A. (2018). Kitsune: An ensemble of autoencoders for online network intrusion detection. In *Proceedings of the Network and Distributed System Security Symposium*. Internet Society. <https://doi.org/10.14722/ndss.2018.23204>
- Montavon, G., Samek, W., & Muller, K. R. (2018). Methods for interpreting and understanding deep neural networks. *Digital Signal Processing*, 73, 1-15. <https://doi.org/10.1016/j.dsp.2017.10.011>
- Moustafa, N., & Slay, J. (2015). UNSW-NB15: A comprehensive data set for network intrusion detection systems. In *2015 Military Communications and Information Systems Conference* (pp. 1-6). IEEE. <https://doi.org/10.1109/MilCIS.2015.7348942>
- Naseer, S., Saleem, Y., Khalid, S., Bashir, M. K., Han, J., Iqbal, M. M., & Han, K. (2018). Enhanced network anomaly detection based on deep neural networks. *IEEE Access*, 6, 48231-48246. <https://doi.org/10.1109/ACCESS.2018.2863036>
- Newman, M. E. J. (2003). The structure and function of complex networks. *SIAM Review*, 45(2), 167-256. <https://doi.org/10.1137/S003614450342480>
- Newman, M. E. J. (2006). Modularity and community structure in networks. *Proceedings of the National Academy of Sciences*, 103(23), 8577-8582. <https://doi.org/10.1073/pnas.0601602103>
- Newman, M. E. J., & Girvan, M. (2004). Finding and evaluating community structure in networks. *Physical Review E*, 69(2), 026113. <https://doi.org/10.1103/PhysRevE.69.026113>
- Patcha, A., & Park, J. M. (2007). An overview of anomaly detection techniques: Existing solutions and latest technological trends.

- Computer Networks, 51(12), 3448-3470. <https://doi.org/10.1016/j.comnet.2007.02.001>
- Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, 1(1), 81-106. <https://doi.org/10.1007/BF00116251>
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 1135-1144). Association for Computing Machinery. <https://doi.org/10.1145/2939672.2939778>
- Ring, M., Wunderlich, S., Scheuring, D., Landes, D., & Hotho, A. (2019). A survey of network-based intrusion detection data sets. *Computers & Security*, 86, 147-167. <https://doi.org/10.1016/j.cose.2019.06.005>
- Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5), 206-215. <https://doi.org/10.1038/s42256-019-0048-x>
- Samek, W., Wiegand, T., & Müller, K. R. (2017). Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *IEEE Signal Processing Magazine*, 34(6), 50-68. <https://doi.org/10.1109/MSP.2017.2765204>
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2020). Grad-CAM: Visual explanations from deep networks via gradient-based localization. *International Journal of Computer Vision*, 128(2), 336-359. <https://doi.org/10.1007/s11263-019-01228-7>
- Sharafaldin, I., Lashkari, A. H., & Ghorbani, A. A. (2018). Toward generating a new intrusion detection dataset and intrusion traffic characterization. In *Proceedings of the 4th International Conference on Information Systems Security and Privacy* (pp. 108-116). SciTePress. <https://doi.org/10.5220/0006639801080116>
- Shone, N., Ngoc, T. N., Phai, V. D., & Shi, Q. (2018). A deep learning approach to network intrusion detection. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2(1), 41-50. <https://doi.org/10.1109/TETCI.2017.2772792>
- Sicari, S., Rizzardi, A., Grieco, L. A., & Coen-Porisini, A. (2015). Security, privacy and trust in Internet of Things: The road ahead. *Computer Networks*, 76, 146-164. <https://doi.org/10.1016/j.comnet.2014.11.008>
- Sisinni, E., Saifullah, A., Han, S., Jennehag, U., & Gidlund, M. (2018). Industrial Internet of Things: Challenges, opportunities, and directions. *IEEE Transactions on Industrial Informatics*, 14(11), 4724-4734. <https://doi.org/10.1109/TII.2018.2852491>
- Sommer, R., & Paxson, V. (2010). Outside the closed world: On using machine learning for network intrusion detection. In *2010 IEEE Symposium on Security and Privacy* (pp. 305-316). IEEE. <https://doi.org/10.1109/SP.2010.25>
- Tjoa, E., & Guan, C. (2021). A survey on explainable artificial intelligence (XAI): Toward medical XAI. *IEEE Transactions on Neural Networks and Learning Systems*, 32(11), 4793-4813. <https://doi.org/10.1109/TNNLS.2020.3027314>
- Vinayakumar, R., Alazab, M., Soman, K. P., Poornachandran, P., Al-Nemrat, A., & Venkatraman, S. (2019). Deep learning approach for intelligent intrusion detection system. *IEEE Access*, 7, 41525-41550. <https://doi.org/10.1109/ACCESS.2019.2895334>
- Watts, D. J., & Strogatz, S. H. (1998). Collective dynamics of 'small-world' networks. *Nature*, 393(6684), 440-442. <https://doi.org/10.1038/30918>
- Xu, L. D., He, W., & Li, S. (2014). Internet of Things in industries: A survey. *IEEE Transactions on Industrial Informatics*, 10(4), 2233-2243. <https://doi.org/10.1109/TII.2014.2300753>
- Xu, L. D., Lu, Y., & Li, L. (2021). Embedding blockchain technology into IoT for security: A survey. *IEEE Internet of Things Journal*, 8(13), 10452-10473. <https://doi.org/10.1109/JIOT.2021.3060508>
- Xu, L. D., Xu, E. L., & Li, L. (2018). Industry 4.0: State of the art and future trends. *International Journal of Production Research*, 56(8), 2941-2962. <https://doi.org/10.1080/00207543.2018.1444806>
- Xu, X., Zhang, J., & Small, M. (2008). Superfamily phenomena and motifs of networks induced from time series. *Proceedings of the National Academy of Sciences*, 105(50), 19601-19605. <https://doi.org/10.1073/pnas.0806082105>
- Zargar, S. T., Joshi, J., & Tipper, D. (2013). A survey of defense mechanisms against distributed denial of service (DDoS) flooding attacks. *IEEE Communications Surveys & Tutorials*, 15(4), 2046-2069. <https://doi.org/10.1109/SURV.2013.031413.00127>
- Zarpelao, B. B., Miani, R. S., Kawakani, C. T., & de Alvarenga, S. C. (2017). A survey of intrusion detection in Internet of Things. *Journal of Network and Computer Applications*, 84, 25-37. <https://doi.org/10.1016/j.jnca.2017.02.009>
- Zhang, C., & Lu, Y. (2021). Study on artificial intelligence: The state of the art and future prospects. *Journal of Industrial Information Integration*, 23, 100224. <https://doi.org/10.1016/j.jii.2021.100224>