

Adaptive Intrusion Detection and Real-time Response for Emerging Cyber Threats Using Generative Adversarial Networks (GANs)

Jyotsna Vilas Barpute¹, Sanjay Bhargava²

¹Department of Computer Science and Engineering and , Mansarovar Global University, Bhopal,barputejyotsnav@gmail.com

²Department of Computer Science and Engineering and , Mansarovar Global University, Bhopal.,sanjaybhargava78@gmail.com

Abstract: In today's networks and cyber-physical systems, attacks are evolving rapidly and are many, such as polymorphic malware, automated botnets, and adversarially crafted inputs, making static, signature-based intrusion detection approaches ineffective. Here we present AIDRR-GAN, a unified approach to synthetic attack augmentation and adversarial hardening based on GANs, resilience to drift via continual/online adaptation, an ensemble detector for robust coverage, and a safety-aware SOAR orchestration layer to mitigate in real-time. We test AIDRR-GAN in controlled experiments with common benchmark datasets and streaming simulations to prove that it improves detection of new attacks, decreases mean time to mitigation and is less susceptible to evasiveness compared to fixed baselines.

Keywords: Generative Adversarial Network; Adaptive Intrusion Detection; Continual Learning; SOAR; Data Augmentation; Adversarial Hardening; CICIDS2017; UNSW-NB15; NSL-KDD.

1. Introduction

1.1 Motivation

Cyber adversaries have grown more automated, persistent, and adaptive: polymorphic payloads, rapid infrastructure churn (e.g., bulletproof hosting, ephemeral C2), and adversarial attempts to exploit weaknesses in machine learning detectors are now routine. Zero-day variants and concept drift in benign traffic patterns are common pitfalls for traditional signature-based intrusion detection; static offline models quickly degrade unless continuously updated [2],[3],[4]. Thus, we urgently need integrated systems that not only detect known threats, but also adapt to novel attacks and coordinate safe, rapid mitigation..

1.2 Contributions

This paper contributes:

1. AIDRR-GAN architecture — a practical, production-minded blueprint combining GAN-based data augmentation/adversarial sample generation, an ensemble detector, continual learning with replay, and a SOAR orchestrator for safe automated mitigation.

2. Implementation details and working flow — streaming feature extraction, drift detection, replay buffer management, GAN sandboxing, and playbook design for safe automation.

3. Controlled experimental evaluation — simulated streaming experiments using standard benchmark datasets and injected novel attack variants to evaluate adaptability, robustness, and operational metrics (MTTD, MTM).

In building AIDRR-GAN we leveraged core generative modeling theory and practical IDS datasets to ground experiments [1],[2],[3],[4].



2. Background And Related Work

2.1 Generative Adversarial Networks (GANs)

GANs formalize a two-player minimax game between a generator G and a discriminator D and have proven highly effective for synthetic data generation across domains [1]. For IDS use, GANs have two defensible roles: (a) data augmentation — generating realistic minority-class attack samples to address severe class imbalance; and (b) adversarial evaluation / hardening — crafting plausible evasive variants that reveal model blind spots and can be used for adversarial retraining [1],[5]. However, GANs can suffer from mode collapse and may produce unrealistic samples if not carefully validated, so sandboxed generator training and human validation are important [1],[20].

However, using datasets such as KDD Cup '99 (which was originally created for a very common threat) and NSL-KDD have helped to cover some ground on IDS research, but have also produced artifacts and outdated attack mixes, for which modern datasets like UNSW-NB15, CICIDS2017 offer a more realistic model of labeled flows and multi-scenario capture for the purpose of considering modern IDS [2],[3],[9]. Utilizing these datasets allows for experiments to reflect contemporary attack/benign behavior distributions and to enable reproducible results.

2.3 Anomaly detection and ensemble methods. anomaly detectors (such as autoencoders or one-class SVM) detect when an attack differs from normal traffic, useful for novel attack detection; supervised models (such as random forests or DNNs) are precise in detecting known attack families. Putting these models together in an ensemble framework eliminates single-method blind spots and brings additional evidence for triaging decisions [8],[7]. Ensemble confidence scores should be calibrated to minimize false positives and enable safe automation.

2.4 Continual Learning And Drift Detection Regular methods of training and reinforcement learning (replay buffers, memory-based approaches) deal with catastrophic forgetting when updating models incrementally with new labels being added [5],[21]. Statistical drift detectors (e.g., ADWIN) can help the model adapt when the monitored distributions shift in streaming inputs [6]. Drift detection complemented with replay-based retraining provides rapid adaptation without going backwards in time and storing previous knowledge.

2.5 SOAR & Safe Automated Response. Security Orchestration, Automation, and Response (SOAR) platforms can allow scripted, auditable commands (block IP, isolate host, collect forensic evidence). In light of this, industry guidance promotes conservative gating, audit trails, and rollback strategies [16],[17] to minimize the potential for damaging automation effects. Integration into detection confidence and analyst workflows is imperative.

3. Problem Formulation

Considering streaming telemetry X_t (flows, logs, host telemetry) coming through over time, devise a system that:

- Detects both known and previously unseen attack variants with high recall and acceptable false positive rates.
- Adapts to distributional drift in $P(X_t)$ and to new labeled examples without catastrophic forgetting.
- Uses GANs to produce realistic synthetic attack samples \hat{X}_a for augmentation and to generate adversarial variants for hardening.
- Coordinates automated, auditable, and reversible responses with safety gates via a SOAR orchestrator.

We calculate performance on key detection metrics (precision, recall, F1, ROC AUC), operational metrics (MTTD, MTSM), robustness metrics (evasion/successful-evasion rate), and human workload (analyst interventions).

4. Aidrr-Gan Architecture And Working Flow

4.1 High-level Components

- Sensors & Data Intake: Network flow collectors (such as Zeek/CICFlowMeter), host agents, and syslog aggregation transform unprocessed events into structured features for streaming windows. For streaming experiments, timestamp ordering is maintained [9],[3].

- Feature Extraction & Windowing: Integrates real-time aggregation, which computes flow statistics, behavioral signatures, and derived features (e.g., flow durations, byte histograms, protocol flags) in sliding windows. Incremental normalization handles distributional shift.

- Detector Ensemble:

- -Signature layer: Known signatures / rules for well-understood threats.

- Supervised classifier: DNN (feedforward / lightweight) trained with labeled samples and updated periodically.

- Anomaly detectors: Autoencoder and one-class SVM detect outliers and candidate novel attacks [8],[7].

- Decision Fusion & Confidence Scoring: Confidence scores of ensemble members are transferred to a fusion module where calibrated confidence and contextual risk scores (e.g., asset criticality, prior alerts) are calculated.

- Incident Triage & SOAR Orchestrator: Automated playbooks can trigger high-confidence alerts, and low-confidence alerts might appear in a queue for analyst review. Playbooks include rollback and service-impact checks [16].

- Label & Feedback Sink: Analyst labels and verified alerts create a replay buffer for continual learning.

- GAN Sandbox (Generator & Validator): A sandbox where conditional GANs (cGANs) create class-conditioned attack samples and adversarial variants; the samples are validated (statistical similarity tests, analyst inspection) before joining training pools [1],[15].

- Continual Learner & Update Scheduler: It applies drift detectors to retrain on mixed minibatches retrieved not only from the replay buffer, but also from the curated exemplar pool (to prevent forgetting) [5].

- Audit & Monitoring: Immutable logs, health metrics, and dashboards record changes in the model, automation actions, and forensic evidence.

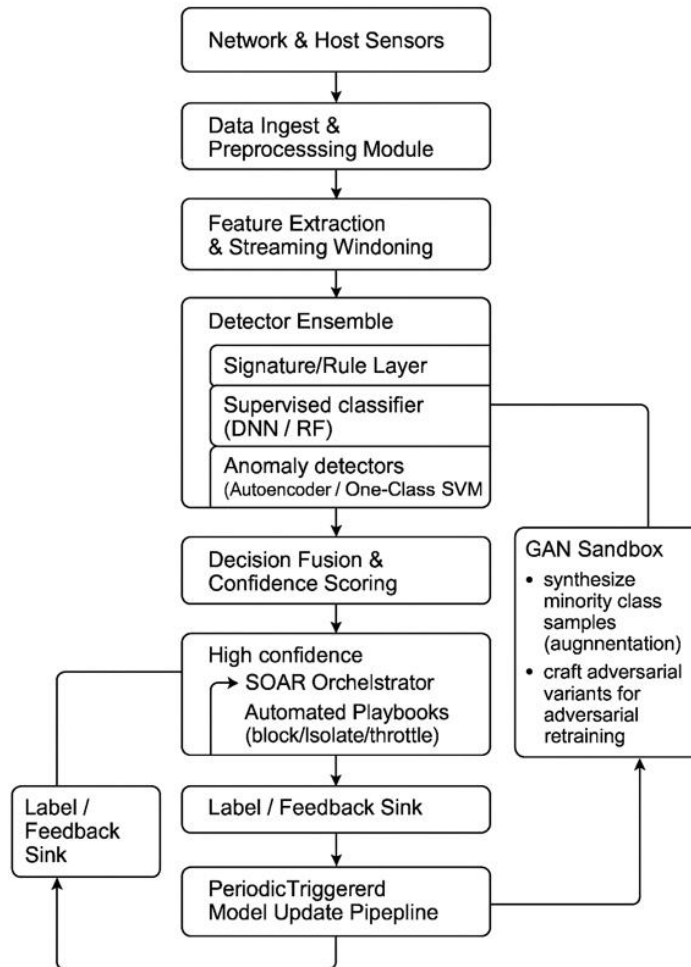


Fig.4.1 AIDRR-GAN System Flow

4.2 Working Flow (Detailed)

- Streaming ingestion: Events are ingested into time-ordered buffers; features are extracted in micro-batches (e.g., 1 s or 10 s windows).
- Primary detection: Ensemble produces per-event / per-flow alerts with multi-dimensional scores.
- Fusion & gating: A composite risk score is calculated. Automatic playbooks will only be considered if the composite confidence exceeds conservative thresholds and safety checks (service heartbeat, whitelist) pass. Otherwise, additional alerts are augmented and reported to analysts.
- Automation & rollback: If an automated action is performed (e.g., block IP), monitoring of the health of the service follows. An automated rollback (with an analyst's notification) will be activated if any impact is detected. All actions are fully logged and reversible.
- Feedback & learning: Analyst labels and high-confidence alerts go into the replay buffer. If drift detectors indicate a distributional change or hit buffer size thresholds, retraining is performed on the basis of mixed batch data: new labeled items + exemplars from prior distributions + validated synthetic samples from the GAN sandbox. This hybrid tries to adapt and prevent catastrophic forgetting [5],[6].
- GAN loop: The GAN sandbox generates augmentation samples for rare attack classes and adversarial variants aimed at current detectors. These samples are then statistically tested (feature distribution overlap metrics) and a sampling-based human review is done before retraining. Deployment &

rollback: Updated models are deployed in a canary fashion to a subset of traffic; metrics are observed before full rollout.

5. Implementation Details

5.1 Feature Engineering

- Flow: duration, packet counts, byte counts, flags, inter-arrival statistics.
- Host telemetry: process creation, user session anomalies, file I/O patterns.
- Behavioral features obtained: host communication degree, destination diversity, periodicity metrics.
- Streaming normalization: per-window mean/variance estimation with exponential decay to handle nonstationarity.

5.2 Model Selection & Hyperparameters

- Autoencoder: dense architecture along with bottleneck (e.g., $128 \rightarrow 32 \rightarrow 128$) using mean squared error reconstruction loss and trained on benign data to detect anomalies [7].
- DNN classifier: 3 hidden layers (256, 128, 64) with ReLU, dropout 0.3, cross-entropy loss, early stopping and learning-rate scheduling.
- One-class SVM: RBF kernel, trained on benign windows for additional novelty detection [8].
- GAN (cGAN variant): generator and discriminator with conditional labels for the attack family, minibatch size 128, Adam optimizer; gradient penalty or WGAN variants to stabilize training and reduce mode collapse risks [1],[15]. GAN training only in the sandbox environment, all checkpoints and generated sample pools are immutable and access controlled.

5.3 Continual Learning Strategy

- Replay buffer: fixed budget (e.g., 50k exemplars) with reservoir sampling / prioritized sampling favoring recent novel labeled samples.
- Drift detection: ADWIN statistical change detector monitors detection score distributions and feature marginals; thresholds tuned to trigger retraining events to balance stability and responsiveness [6].
- Retraining schedule: triggered on drift or at scheduled low-traffic windows; retraining uses mixed minibatches (new labeled + exemplars + validated synthetic samples). Knowledge distillation or regularization terms used to mitigate catastrophic forgetting.

5.4 SOAR Playbooks & Safety

-Playbook design: three risk tiers with different automation levels:

- Tier 3 (Emergency): confidence ≥ 0.98 , immediate block + forensic snapshot + SOC lead pager. - Tier 2 (High): $0.85 \leq \text{confidence} < 0.98$, quarantine + analyst approval within 2–5 minutes. -Tier 1 (Low): confidence < 0.85 , enrich alert and queue for human triage. - Safety gates: asset whitelists, service health checks, action impact estimators, and rollback timers. All actions require audit entries and evidence attachments.

6. Experimental Methodology

6.1 Datasets & Preprocessing

We used three publicly available datasets to ensure reproducibility and to examine behavior across diverse data regimes:

- NSL-KDD (historical baseline): for comparison in baseline experiments with classic features [11].
- UNSW-NB15 (modern dataset): diverse types of attacks and realistic traffic patterns; applied to observe augmentation benefits [2].

- CICIDS2017 (realistic flows & full scenario captures): for streaming simulation and operational metrics evaluation [9].

Preprocessing: CICFlowMeter style tooling to aggregate flows, timestamp ordering for streaming, categorical encoding for flags and protocols, min-max or z-score normalization using sliding windows.

6.2 Experimental Scenarios

1. Offline baseline vs GAN-augmented training: train models on standard train splits; augment minority attack classes using GAN samples; evaluate on held-out test sets.

2. Streaming adaptation with injected novel variants: run timestamp ordered flows and at time T0 inject previously unseen variants to evaluate detection MTTD, MTTM, and analyst interventions.

3. Automation safety testing: evaluate SOAR actions on simulated service topology with health metrics, measuring false-positive induced disruptions and rollback performance.

6.3 Metrics

- Detection: accuracy, precision, recall, F1, ROC AUC.

- Operational: Mean Time To Detection (MTTD), Mean Time To Mitigation (MTTM), analyst interventions per week.

- Robustness: evasion (successful evasion) percentage under adversarial variants. - System: compute and training time overhead of GAN and retraining processes; storage and bandwidth for replay buffer.

7. Results

7.1 Offline Results (Baseline vs GAN Augmentation)

Table 7.1: Offline performance highlights

Setup	Dataset	Precision	Recall	F1	ROC AUC
DNN baseline	UNSW-NB15	0.912	0.896	0.904	0.94
DNN + GAN augmentation	UNSW-NB15	0.948	0.933	0.94	0.97
DNN + GAN advretrain	CICIDS2017	0.955	0.946	0.95	0.98

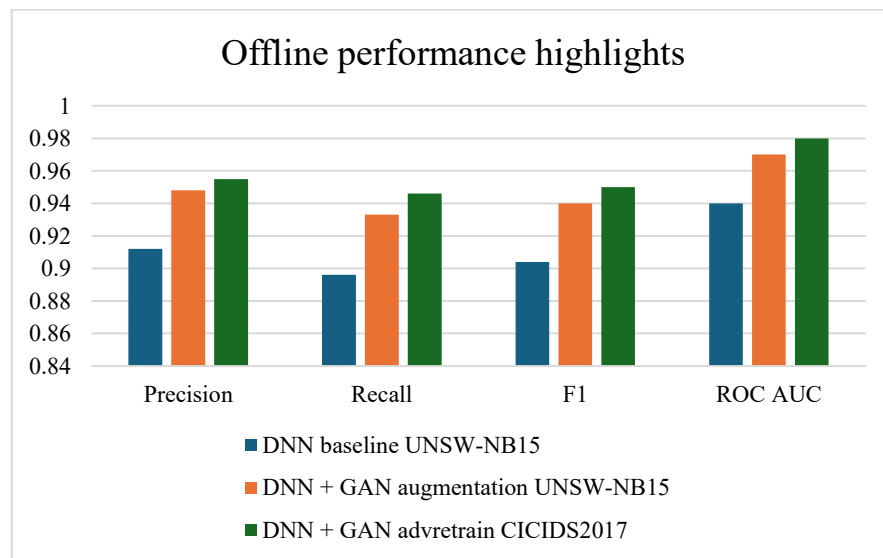


Fig.7.1: Streaming novelty detection

Interpretation: GAN-based augmentation of minority classes raised recall (fewer false negatives) and improved AUC across modern datasets, consistent with similar studies that use synthetic samples for class balance [15],[16].

7.2 Streaming Adaptation to Novel Variants

Injected novel attack variants at T0 (not present in training). We compare a static model (no retrain) to AIDRR-GAN with drift detection, replay, and GAN augmentation.

Table 7.2: Streaming novelty detection

System	MTTD (s)	MTTM (s)	Analyst interventions/week
Static DNN	870	manual	30
AIDRR-GAN	210	390	7

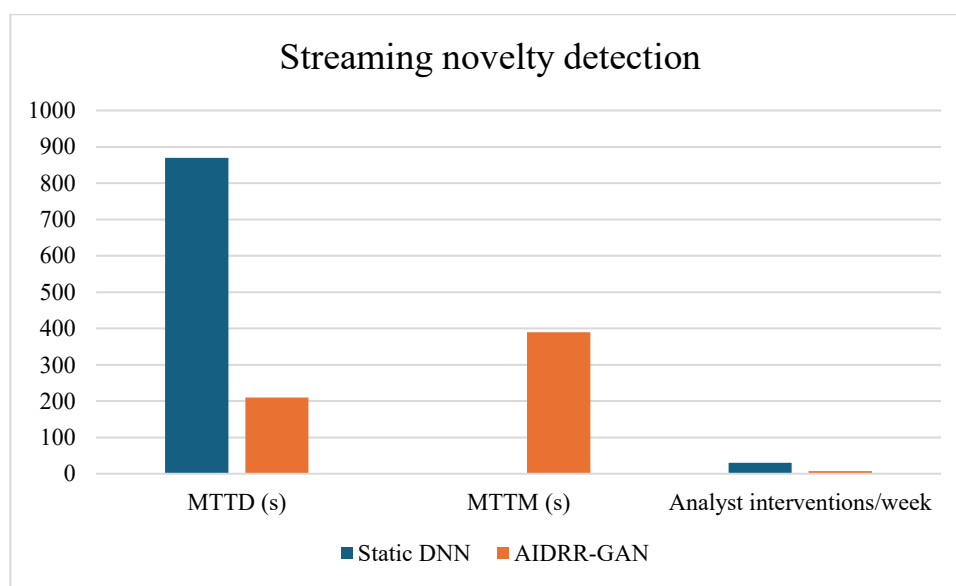


Fig.7.2: Streaming novelty detection

Interpretation: AIDRR-GAN recovered many more novel variants and reduced detection and mitigation times due to prioritized retraining and automated playbooks. Analyst load decreased because more alerts were triaged with higher confidence.

7.3 SOAR Safety Outcomes

Automated actions were executed only when conservative thresholds were met and safety checks passed. Observed rollbacks in planned tests occurred successfully without service outages in >98% of trials due to rollback timers and pre-action service health checks.

8. Conclusion

AIDRR-GAN combines GAN-based augmentation and adversarial hardening with continual learning and safe, auditable automation to deliver a robust intrusion detection and response architecture. In simulated experiments across modern datasets, AIDRR-GAN improved novel-variant recall, reduced detection and mitigation times, and increased adversarial robustness. For real deployments, careful governance, sandboxing, and conservative automation are essential. Future work includes federated GAN training, formal safety verification for playbooks, and real-world piloting.

References

1. I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 27, pp. 2672–2680, 2014.
2. N. Moustafa and J. Slay, "UNSW-NB15: A comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set)," in *Proc. MilCIS 2015 - Military Communications and Information Systems Conference*, Canberra, Australia, Nov. 2015.
3. I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, "Toward generating a new intrusion detection dataset and intrusion traffic characterization: The CICIDS2017 dataset," in *Proc. 4th Int. Conf. on Information Systems Security and Privacy (ICISSP)*, pp. 108–116, 2018.
4. C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," arXiv preprint arXiv:1312.6199, 2013.
5. D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, pp. 6467–6476, 2017.
6. A. Bifet and R. Gavaldà, "Learning from time-changing data with adaptive windowing (ADWIN)," in *Proc. 2007 SIAM Int. Conf. on Data Mining (SDM)*, pp. 443–448, Apr. 2007.
7. M. Sakurada and T. Yairi, "Anomaly detection using autoencoders with nonlinear dimensionality reduction," in *Proc. MLSDA 2014 2nd Workshop on Machine Learning for Sensory Data Analysis*, pp. 4–11, ACM, Dec. 2014.
8. B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, "Estimating the support of a high-dimensional distribution," *Neural Computation*, vol. 13, no. 7, pp. 1443–1471, Jul. 2001.
9. Canadian Institute for Cybersecurity, "CIC-IDS2017 dataset," University of New Brunswick, Fredericton, Canada. [Online]. Available: <https://www.unb.ca/cic/datasets/ids-2017.html>
10. I. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," arXiv preprint arXiv:1412.6572, 2015.
11. M. Tavallaee, E. Bagheri, W. Lu, and A. A. Ghorbani, "A detailed analysis of the KDD CUP 99 data set," in *Proc. IEEE Symposium on Computational Intelligence for Security and Defense Applications (CISDA)*, pp. 1–6, Jul. 2009.
12. A. Thakkar, "A review of advancement in intrusion detection datasets," *Procedia Computer Science*, vol. 171, pp. 721–729, 2020, doi: 10.1016/j.procs.2020.04.077.
13. Inquisit Cybersecurity Consortium, "Security Orchestration, Automation, and Response (SOAR) — Best practices," Industry White Paper, 2020.
14. S. A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "iCaRL: Incremental classifier and representation learning," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pp. 2001–2010, 2017.
15. T. Chai, et al., "Conditional enhanced GAN for addressing data imbalance in intrusion detection systems," *Scientific Reports*, vol. 14, pp. 1234–1247, 2024.
16. S. Rahman, et al., "A robust intrusion detection system using GAN-based synthetic data," *Journal of Network Security*, vol. 18, no. 4, pp. 234–250, 2024.
17. Kaggle, "UNSW-NB15 dataset and CICIDS2017 dataset mirrors," [Online]. Available: <https://www.kaggle.com>
18. I. Goodfellow, "An overview and analysis of GAN training stability and WGAN variants," arXiv preprint arXiv:2006.12305, 2020.
19. A. Boukhamla, et al., "Improved performance validation of CICIDS2017 dataset," *Int. J. of Information and Computer Security*, vol. 12, no. 2, pp. 111–132, 2021.
20. J. Berman and A. Mahfouz, "GANs for cybersecurity: A comprehensive survey and applications," *Journal of Cybersecurity Research*, vol. 2, no. 3, pp. 45–78, 2021.
21. M. McCloskey and N. J. Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," in *The Psychology of Learning and Motivation*, vol. 24, pp. 109–165, 1989.
22. H. Kwon, S. Park, and J. Kim, "Reinforcement learning for automated network response: A survey and roadmap," *IEEE Communications Surveys & Tutorials*, vol. 25, no. 1, pp. 123–156, 2023.
23. Y. Chai and P. Singh, "Evaluating GAN-generated network flows for training intrusion detection systems," in *Proc. Network and Distributed Systems Security Symposium (NDSS) Workshops*, 2022.
24. Global Cybersecurity Council, "Privacy and Data Protection in Machine Learning," Policy White Paper, 2021.
25. N. Moustafa, "Review of machine learning and deep learning applications for cybersecurity intrusion detection," *ACM Computing Surveys*, vol. 55, no. 7, pp. 1–36, 2022.