

# LATENT DIFFUSION FOR INTERNET OF THINGS ATTACK DATA GENERATION IN INTRUSION DETECTION \*

A PREPRINT

✉ Estela Sánchez-Carballo<sup>†1</sup>, ✉ Francisco M. Melgarejo-Meseguer<sup>1</sup>, and ✉ José Luis Rojo-Álvarez<sup>1,2</sup>

<sup>1</sup>Departamento de Teoría de la Señal y Comunicaciones y Sistemas Telemáticos y Computación,  
Universidad Rey Juan Carlos, Madrid, Spain.

<sup>2</sup>Data, Networks and Cybersecurity Research Institute, Universidad Rey Juan Carlos, Madrid, Spain.

## ABSTRACT

Intrusion Detection Systems (IDSs) are a key component for protecting Internet of Things (IoT) environments. However, in Machine Learning-based (ML-based) IDSs, performance is often degraded by the strong class imbalance between benign and attack traffic. Although data augmentation has been widely explored to mitigate this issue, existing approaches typically rely on simple oversampling techniques or generative models that struggle to simultaneously achieve high sample fidelity, diversity, and computational efficiency. To address these limitations, we propose the use of a Latent Diffusion Model (LDM) for attack data augmentation in IoT intrusion detection and provide a comprehensive comparison against state-of-the-art baselines. Experiments were conducted on three representative IoT attack types, specifically Distributed Denial-of-Service (DDoS), Mirai, and Man-in-the-Middle, evaluating both downstream IDS performance and intrinsic generative quality using distributional, dependency-based, and diversity metrics. Results show that balancing the training data with LDM-generated samples substantially improves IDS performance, achieving F1-scores of up to 0.99 for DDoS and Mirai attacks and consistently outperforming competing methods. Additionally, quantitative and qualitative analyses demonstrate that LDMs effectively preserve feature dependencies while generating diverse samples and reduce sampling time by approximately 25% compared to diffusion models operating directly in data space. These findings highlight latent diffusion as an effective and scalable solution for synthetic IoT attack data generation, substantially mitigating the impact of class imbalance in ML-based IDSs for IoT scenarios.

**Keywords** Data Augmentation · Diffusion Models · Internet of Things · Intrusion Detection Systems · Latent Spaces · Tabular Data Synthesis.

## 1 Introduction

The rapid proliferation of Internet of Things (IoT) devices has transformed several domains such as healthcare, smart cities, industrial automation, and intelligent transportation systems [Perwej et al., 2019, Lohiya and Thakkar, 2021]. However, this increase in smart environments has also significantly increased the attack rate of modern networks, making IoT ecosystems particularly vulnerable to cyber-threats [Shafiq et al., 2022]. As a result, Intrusion Detection Systems (IDSs) have become fundamental for protecting IoT infrastructures [Khraisat et al., 2019]. In recent years, IDSs based on Machine Learning (ML) and Deep Learning (DL) have demonstrated strong performance, as they are capable of learning complex patterns from network traffic data [Shenfield et al., 2018, Melgarejo-Meseguer et al., 2021, Huang et al., 2025]. However, their performance highly depends on the availability of large, representative, and well-balanced datasets. A key challenge in this context is that IoT environments are exposed to a wide variety of attack families, ranging from high-rate Distributed Denial-of-Service (DDoS) and botnet-driven flooding attacks to low-rate and stealthy Man-in-the-Middle (MitM) attacks, each inducing different traffic patterns and detection challenges [Kaur et al., 2023]. Additionally, IoT attack datasets suffer from several limitations. First, collecting real-world attack traffic is costly, time-consuming, and often limited by privacy and regulatory concerns [Alex et al., 2023, Dritsas and Trigka, 2025]. Second, these datasets are typically highly imbalanced, with not common but critical attack types being severely underrepresented

\*This work has been submitted to the IEEE for possible publication

<sup>†</sup>estela.sanchezc@urjc.es

compared to benign traffic [Ma et al., 2025, Mallidi and Ramisetty, 2025]. Third, IoT traffic data are heterogeneous by nature, combining tabular continuous (e.g., packet sizes, timing statistics) and categorical (e.g., protocol flags, connection states) features [Kaur et al., 2023, Rahman et al., 2025]. These characteristics make training robust IDSs more difficult and motivate the use of synthetic data generation techniques as a complementary strategy to data collection.

To address data scarcity and imbalance in IoT attack datasets, a wide range of synthetic data generation techniques have been proposed in the literature. Classical oversampling methods, such as the Synthetic Minority Oversampling Technique (SMOTE), generate new samples by interpolating between minority-class instances [Chawla et al., 2002]. While computationally efficient, these approaches operate locally in feature space and generally maintain pairwise feature dependencies, but tend to generate new samples that are highly similar to the original data, resulting in limited diversity [Azhar et al., 2022]. More recently, deep generative models have been adopted for tabular data synthesis in cybersecurity, including Variational Autoencoders (VAEs) [Zhang et al., 2023] and Generative Adversarial Networks (GANs) [Wang et al., 2017]. These models are capable of learning global data distributions and capturing higher-order correlations, but they also have limitations. VAEs can generate over-smoothed samples due to the imposed latent regularization, and GANs are known to suffer from training instability, mode collapse, and sensitivity to hyperparameter tuning issues [Wang et al., 2017, 2024, Li et al., 2025a]. Diffusion Models (DMs) have recently emerged as a powerful generative modeling approach, achieving state-of-the-art results in multiple domains by progressively denoising samples from pure noise [Ho et al., 2020]. Their strong theoretical foundations and training stability make them attractive for tabular data generation, however, their direct application to tabular IoT data remains challenging. Operating directly in the feature space requires careful handling of heterogeneous variables, leads to slow sampling due to the high dimensionality of IoT traffic data and the large number of diffusion steps involved, and can struggle with features that exhibit different statistical scales [Wang et al., 2024, Li et al., 2025b].

In this work, we propose leveraging the advantages of DMs over existing generative approaches while addressing their known limitations in the context of IoT traffic data. To this purpose, we introduce a latent diffusion framework in which a mixed-variable Autoencoder (AE) is first employed to encode heterogeneous network traffic features into a compact and normalized latent space, and a DM is then trained to capture the distribution of these latent representations. Performing the diffusion process in the latent domain results in improved training stability and significantly reduced computational complexity, while providing a more uniform representation across features with different statistical properties. At the same time, the AE decoder preserves the semantic structure of the original data, enabling the generation of realistic and coherent synthetic IoT traffic samples. Latent DMs (LDMs) were originally introduced for image synthesis [Rombach et al., 2022] and have recently shown promising results in tabular data generation [Sattarov et al., 2023, Villaizán-Vallelado et al., 2025, Zhu et al., 2025]. A key contribution of this work is to transfer this latent diffusion paradigm to the IoT traffic domain and to demonstrate its effectiveness for synthetic IoT attack data generation in intrusion detection.

## 2 Related Work

Class imbalance is a well-known challenge in IoT IDSs, where benign traffic typically dominates and many attack types are severely underrepresented [Mallidi and Ramisetty, 2025]. Classical oversampling techniques, such as SMOTE, have been widely used to address this issue. For example, in [Soe et al., 2019] and [Dash et al., 2020], SMOTE was applied to resolve extreme data imbalance in different IoT attack datasets prior to training DL-based IDSs. However, classical oversampling techniques have been reported to provide limited performance gains in highly imbalanced datasets with mixed continuous and categorical features [Wang et al., 2024], which has motivated the use of deep generative models capable of improving our capabilities for learning global data distributions.

VAEs emerged as an alternative to classical oversampling techniques for synthetic tabular data generation. They generate new samples by learning probabilistic latent representations of the data and reconstructing them through a decoder network [Pinheiro-Cinelli et al., 2021]. Several studies have explored the use of VAEs for generating synthetic IoT traffic data to mitigate class imbalance in intrusion detection. For instance, in [Liu et al., 2022], two different VAE variants were evaluated to balance an IoT traffic dataset with the aim of improving IDS performance. Similarly, in [Zhang and Liu, 2022], a VAE-based approach was proposed to generate synthetic IoT attack samples, also reporting improved IDS performance after data augmentation. However, VAEs are known to produce over-smoothed samples due to latent space regularization, which can limit their effectiveness for synthesizing highly imbalanced and complex tabular IoT attack data [Pinheiro-Cinelli et al., 2021].

GANs represent another DL-based alternative to classical oversampling techniques for synthetic tabular data generation by learning the underlying data distribution through adversarial training between a generator and a discriminator [Wang et al., 2017]. In the context of IoT security, GANs have been employed to mitigate class imbalance in intrusion detection datasets. For instance, in [Park et al., 2023], authors observed that ML models struggle to learn meaningful attack patterns under severe class imbalance and used a GAN to generate synthetic IoT traffic samples. Similarly, in [Rahman et al., 2024], authors utilized a GAN for synthesizing IoT attack data, reporting significant improvements in IDS performance. However, GANs are known to suffer from training

instability, mode collapse, and sensitivity to architectural and optimization choices, which can limit their robustness [Wang et al., 2017].

In recent years, DMs have emerged as a powerful class of generative models, addressing several limitations of VAEs and GANs [Wang et al., 2024]. DMs generate new data by learning a denoising process that progressively transforms samples from pure noise into realizations drawn from the target data distribution [Ho et al., 2020]. In the context of IoT security, DMs have recently been explored for synthetic attack data generation. For instance, in [Camerota et al., 2024], a DM was proposed to generate synthetic IoT attack samples, showing that malware detection performance improves when the generated data are incorporated into training. Similarly, in [Wang et al., 2025], a diffusion-transformer framework in which a DM was used to balance the dataset, followed by a transformer-based classifier was introduced to enhance IDS performance. However, DMs are often computationally expensive and sensitive to feature heterogeneity, which can limit their practicality for large-scale and high-dimensional IoT attack datasets [Li et al., 2025b].

A common conclusion across the reviewed literature is that synthetic data generation is essential for IoT intrusion detection, as attack classes are typically underrepresented, leading to poor IDS performance that has been shown to substantially improve when synthetic samples are incorporated. Among the generative approaches studied, DMs represent the most advanced techniques and have demonstrated the ability to generate high-quality synthetic IoT attack data [Camerota et al., 2024, Wang et al., 2025]. However, their high computational cost and their sensitivity to mixed-type features and variables with different statistical scales leave room for further improvement [Wang et al., 2024]. To address these limitations, LDMs [Rombach et al., 2022] have been proposed, in which the diffusion process is learned in a low-dimensional space obtained through an AE. While LDMs have been successfully applied to tabular data synthesis in other domains, such as financial data [Sattarov et al., 2023] and sensitive personal data including demographic, housing, and health-related information [Zhu et al., 2025], their application to IoT network traffic data has not yet been explored. In these domains, LDMs have been shown to perform competitively with state-of-the-art generative models, motivating their investigation for IoT attack data generation.

### 3 Methodology

In this section, the proposed method for generating synthetic IoT traffic data is detailed. First, the problem under consideration is introduced. Second, the design of an AE capable of handling mixed-type IoT tabular data is presented, followed by a description of the latent diffusion process. The latent sampling and decoding procedures are then explained, and finally, an analysis of the computational complexity of the proposed framework is provided.

#### 3.1 Problem Formulation

Let  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$  denote an IoT network traffic dataset. Each input sample  $\mathbf{x}_i$  is composed of continuous and categorical variables,  $\mathbf{x}_i = (\mathbf{x}_i^{(c)}, \mathbf{x}_i^{(b)})$ , where  $\mathbf{x}_i^{(c)} \in \mathbb{R}^{d_c}$  consists of  $d_c$  continuous features (e.g. packet sizes, inter-arrival times, and statistical summaries), and  $\mathbf{x}_i^{(b)} \in \{0, 1, \dots, N_b\}^{d_b}$  consists of  $d_b$  binary or categorical features (e.g., protocol flags and connection states). The corresponding label is denoted by  $y_i \in \mathcal{Y}$ , indicating benign traffic or a specific attack type. In this work, we focus on class-specific synthetic data generation.

Given a real dataset  $\mathcal{D}_{real}$ , the goal of this study is to develop a generative model capable of producing synthetic samples  $\hat{\mathbf{x}} \sim p_{\theta}(\mathbf{x}|y)$  such that the resulting augmented dataset  $\mathcal{D}_{aug} = \mathcal{D}_{real} \cup \mathcal{D}_{syn}$  preserves the statistical structure, feature dependencies, and semantic consistency of real IoT traffic. For this purpose, the generated samples should follow the joint distribution of original features, be consistent with the target class, and improve data coverage while avoiding unrealistic samples or points that lie outside the true data manifold.

Directly modeling the joint distribution  $p(\mathbf{x})$  in the original feature space is challenging due to the heterogeneous nature of IoT traffic data, differences in statistical scales across features, and the high dimensionality of the data [Kotelnikov et al., 2023, Wang et al., 2024]. In addition, DMs operating in data space require iterative denoising procedures with high computational cost and sensitivity to feature normalization [Nichol and Dhariwal, 2021, Karras et al., 2022]. To address these challenges, we propose performing the sample generation process in a learned latent space, following the latent diffusion paradigm and adapting it to the specific characteristics of mixed-type IoT traffic data. Original data dimensions are reduced by an AE encoder, and a DM is then trained to approximate the latent distribution, from which new latent samples can be drawn and decoded back to the original feature space via an AE decoder. Under this formulation, the synthetic data generation task reduces to learning a diffusion-based generative process in latent space, followed by a deterministic decoding step.

#### 3.2 Autoencoder for Mixed-Type Tabular Data

To enable the DM to operate in a lower-dimensional space, the original data must first be reduced in dimensionality. For this purpose, we employ an AE. The role of the AE is to learn a compact and continuous latent representation of the original mixed-type

Table 1: Summary of the main characteristics of the attacks considered.

Attack	Category	Traffic Volume	Visibility	Primary Effect
<i>DDoS-ICMP_Fragmentation</i>	Volumetric flooding	Very high	Highly visible	Resource exhaustion
<i>Mirai-greip_flood</i>	Botnet flooding	High, distributed	Moderately visible	Bandwidth/CPU exhaustion
<i>MITM-ArpSpoofing</i>	Protocol manipulation	Low	Stealthy	Traffic interception

IoT traffic data that is suitable for diffusion-based generative modeling. The AE encoder  $f_\phi(\cdot)$  maps the input vector  $\mathbf{x}_i$  into a latent representation  $\mathbf{z}_i = f_\phi(\mathbf{x}_i) \in \mathbb{R}^{d_z}$ , where  $d_z \ll d_c + d_b$ . The encoder is implemented as a Multilayer Perceptron (MLP) with fully connected layers and nonlinear activation functions, enabling the model to capture complex interactions among heterogeneous features. To ensure that the latent space is well conditioned for the diffusion process, the encoder outputs are normalized to have zero mean and unit variance across latent dimensions, computed over the training set. This normalization mitigates scale imbalances across latent variables and improves stability during subsequent diffusion training. The decoder  $g_\phi(\cdot)$  maps latent representations back to the original feature space  $\hat{\mathbf{x}}_i = g_\phi(\mathbf{z}_i)$ , and is designed to reconstruct continuous and discrete features using appropriate output parameterizations. Continuous features are reconstructed using real-valued outputs, while binary and categorical features are reconstructed through sigmoid and softmax activations, respectively. This design allows the AE to explicitly account for the mixed-type nature of IoT traffic data during reconstruction.

The AE is trained to minimize a reconstruction loss that reflects the different statistical properties of continuous and discrete variables. Specifically, the loss function is defined as

$$\mathcal{L}_{AE} = \mathcal{L}_c(\mathbf{x}^{(c)}, \hat{\mathbf{x}}^{(c)}) + \mathcal{L}_b(\mathbf{x}^{(b)}, \hat{\mathbf{x}}^{(b)}), \quad (1)$$

where  $\mathcal{L}_c(\cdot)$  denotes a mean squared error (MSE) loss for continuous features, and  $\mathcal{L}_b(\cdot)$  denotes a binary or categorical cross-entropy loss for discrete features. This formulation ensures balanced learning across feature types and prevents dominance by either continuous or discrete variables. The learned latent representation serves as the input space for the DM. By mapping heterogeneous tabular data into a compact, continuous, and normalized latent space, the AE simplifies the learning task of the DM and improves training stability. Additionally, separating representation learning (AE) from generative modeling (DM) allows each component to focus on a well-defined objective, resulting in a more scalable and robust synthetic data generation framework.

### 3.3 Latent Diffusion

Once the mixed-type IoT traffic data have been mapped into a compact latent representation by the AE, a DM is employed to learn the distribution of the latent variables. Let  $\mathbf{z}_0 \in \mathbb{R}^{d_z}$  denote a generic latent representation obtained from the encoder  $f_\phi(\cdot)$ . The objective of the DM is to learn the underlying data distribution  $p(\mathbf{z}_0)$  by gradually corrupting latent samples with noise and then learning to reverse this process. The diffusion process is divided into two stages, i.e., the forward and the reverse (generative) diffusion processes.

The forward diffusion process is defined as a fixed Markov chain that progressively adds Gaussian noise to the latent variables over  $T$  timesteps. At each timestep  $t \in \{1, 2, \dots, T\}$ , the latent variable  $\mathbf{z}_t$  is obtained from  $\mathbf{z}_{t-1}$  as

$$\mathbf{z}_t = \sqrt{\alpha_t} \mathbf{z}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_t, \quad \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (2)$$

where  $\{\alpha_t\}_{t=1}^T$  is a predefined noise schedule with  $0 < \alpha_t < 1$ . Defining  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ , the forward process admits a closed-form expression that allows sampling  $\mathbf{z}_t$  directly from  $\mathbf{z}_0$ , as follows,

$$\mathbf{z}_t = \sqrt{\bar{\alpha}_t} \mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (3)$$

As  $t$  increases, the latent variables become increasingly dominated by noise, and for sufficiently large  $T$ ,  $\mathbf{z}_T$  approaches a Gaussian distribution.

The reverse diffusion process, also referred to as the generative process, corresponds to learning the reverse Markov chain that progressively removes noise from  $\mathbf{z}_T$  in order to recover samples from  $p(\mathbf{z}_0)$ . This reverse process is parameterized by a neural network  $\epsilon_\theta(\mathbf{z}_t, t)$ , which is trained to predict the noise added at timestep  $t$ . The reverse transition is given by

$$p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t) = \mathcal{N}(\mathbf{z}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{z}_t, t), \boldsymbol{\Sigma}_t), \quad (4)$$

where the mean  $\boldsymbol{\mu}_\theta(\mathbf{z}_t, t)$  is defined as

$$\boldsymbol{\mu}_\theta(\mathbf{z}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{z}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{z}_t, t) \right), \quad (5)$$

and  $\boldsymbol{\Sigma}_t$  is a fixed variance determined by the noise schedule.

The DM is trained by minimizing a denoising objective that encourages accurate noise prediction at each timestep  $t$ . Specifically, the loss function is defined as

$$\mathcal{L}_{DM} = \mathbb{E}_{\mathbf{z}_0, \epsilon, t} \left[ \left\| \epsilon - \epsilon_\theta(\sqrt{\alpha_t} \mathbf{z}_0 + \sqrt{1 - \alpha_t} \epsilon, t) \right\|_2^2 \right], \quad (6)$$

where  $t$  is sampled uniformly from  $\{1, 2, \dots, T\}$ . This formulation allows the DM to be trained efficiently using stochastic gradient descent while leveraging samples drawn directly from the forward process.

### 3.4 Latent Sampling and Decoding

Synthetic IoT traffic samples are generated by sampling from the learned reverse diffusion process and subsequently decoding the resulting latent representations back into the original feature space. A latent variable  $\mathbf{z}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  is first sampled and then iteratively denoised over  $T$  timesteps according to the learned reverse transitions  $\mathbf{z}_{t-1} \sim p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t)$ , with  $t = T, T-1, \dots, 1$ . After completing the reverse diffusion process, a synthetic latent representation  $\hat{\mathbf{z}}_0$  is obtained, which follows the learned latent data distribution  $p_\theta(\mathbf{z}_0)$ . This iterative denoising procedure enables the DM to generate diverse latent samples while preserving the global structure of the latent space learned by the AE. The synthetic latent representation  $\hat{\mathbf{z}}_0$  is then mapped back to the original feature space using the decoder,  $\hat{\mathbf{x}} = g_\phi(\hat{\mathbf{z}}_0)$ , yielding a synthetic sample  $\hat{\mathbf{x}} = (\hat{\mathbf{x}}^{(c)}, \hat{\mathbf{x}}^{(b)})$ , where  $\hat{\mathbf{x}}^{(c)}$  and  $\hat{\mathbf{x}}^{(b)}$  denote the reconstructed continuous and discrete features, respectively.

### 3.5 Computational Complexity Analysis

Let  $N$  denote the number of training samples,  $d = d_c + d_b$  the dimensionality of the original feature space,  $d_z$  the dimensionality of the latent space, and  $T$  the number of diffusion timesteps. The AE is trained once with a computational cost of  $\mathcal{O}(N \cdot C_{AE})$ , where  $C_{AE}$  denotes the cost of a single forward-backward pass through the encoder-decoder network. During inference, each sample is encoded and decoded only once, while diffusion sampling requires multiple sequential steps, making the computational cost of encoding and decoding negligible in comparison. Thus, the dominant computational cost arises from training and sampling the DM. In data-space DMs, this cost scales with the original feature dimension  $d$ , yielding a complexity of  $\mathcal{O}(T \cdot C_{DM}(d))$ . In contrast, the proposed approach operates in latent space, resulting in a reduced complexity of  $\mathcal{O}(T \cdot C_{DM}(d_z))$ . Since  $d_z \ll d$ , this leads to a substantial reduction in both training and sampling costs. Therefore, although the proposed framework introduces an additional AE training step, this cost is incurred only once and is outweighed by the substantial gains in efficiency and scalability achieved by performing diffusion in a latent space.

## 4 Experimental Setup

This section describes the experimental setup to evaluate the proposed approach. First, the dataset used in the experiments and the preprocessing steps applied to the input data are detailed. Then, the hyperparameter optimization strategy employed for all DL models is described. Next, the implementation of the baseline models used for comparison with the LDM is presented. Finally, the architecture and configuration details of the LDM are reported.

### 4.1 Dataset

In this work, the CIIoT2023 dataset released by the Canadian Institute for Cybersecurity, University of New Brunswick [Pinto-Neto et al., 2023], was used. This dataset captures realistic IoT network traffic generated in a controlled environment involving multiple IoT devices and includes both benign traffic and a diverse set of attack scenarios. In this study, the analysis was conducted using pre-extracted flow-level features provided in CSV format, which were designed for ML evaluation. For a detailed description of these features, please refer to [Pinto-Neto et al., 2023].

We focused on three attack types, namely, *DDoS-ICMP\_Fragmentation*, *Mirai-greip\_flood*, and *MITM-ArpSpoofing*. These attacks were selected because they represent different attack families and operate through different mechanisms that induce different traffic patterns in IoT environments. *DDoS-ICMP\_Fragmentation* is a DDoS attack that abuses the Internet Control Message Protocol (ICMP) by transmitting a large volume of fragmented ICMP packets toward a target. The required reassembly of packet fragments significantly increases processing and memory requirements, potentially leading to resource collapse [Sonar and Upadhyay, 2014]. *Mirai-greip\_flood* belongs to the Mirai Botnet family, one of the most well-known malware families targeting IoT devices. In this variant, compromised IoT devices act as bots and generate high-rate flooding traffic toward a target, resulting in bandwidth saturation and Central Processing Unit (CPU) collapse [Abbas et al., 2021]. Finally, *MITM-ArpSpoofing* is a MitM attack that takes advantage of vulnerabilities in the Address Resolution Protocol to intercept or manipulate traffic within a local network. This attack is typically low-rate and secretive [Fereidouni et al., 2025]. Table 1 summarizes the main characteristics of the considered attack types.

Table 2: Attack-specific hyperparameters for each generative model.

Model	Hyperparameter	DDoS	Mirai	MitM
VAE	Latent dimension	5	11	14
	Learning rate ( $\cdot 10^{-3}$ )	0.95	1.72	1.40
	KL weight ( $\cdot 10^{-2}$ )	9.77	1.31	1.01
	Gradient clipping	4.22	4.71	0.06
GAN	Latent noise dimension	13	15	12
	Generator learning rate ( $\cdot 10^{-4}$ )	4.53	1.41	2.03
	Discriminator learning rate ( $\cdot 10^{-4}$ )	4.90	0.98	1.16
	Real-label smoothing	0.80	0.80	0.91
DM	Diffusion steps	200	200	200
	Network depth	5	8	8
	Learning rate ( $\cdot 10^{-4}$ )	3.36	0.61	3.00
LDM	AE latent dimension	13	16	11
	AE learning rate ( $\cdot 10^{-3}$ )	2.22	0.46	2.66
	Diffusion steps	782	461	200
	Diffusion network depth	2	5	5
	Diffusion learning rate ( $\cdot 10^{-4}$ )	2.31	0.89	0.52

## 4.2 Preprocessing

A dataset composed of benign traffic samples and samples corresponding to the three considered attacks was created. As an initial step, constant features in the dataset were removed, and the remaining features were categorized into continuous numerical variables and discrete variables. Outlier detection was applied exclusively to continuous features. To avoid damaging the original data distribution, the outlier removal procedure was designed such that less than 1% of the total data was discarded. Specifically, for each numerical feature, values were first normalized to the  $[0, 1]$  range and sorted. The method then identified gaps between consecutive ordered values, and a predefined threshold determined whether a gap was considered sufficiently large or not. Based on this criterion, lower and upper bounds were computed for each feature. A data point was classified as an outlier if at least one of its numerical feature values lied outside the corresponding bounds. After outlier removal, attack samples were separated from benign traffic, as only attack instances were used to train the generative models. The attack data were then split into training (80%), validation (10%) and test (10%) subsets. The training and validation subsets were used for training and validation of the generative models, respectively, while the test set was exclusively used to evaluate the generated synthetic samples. From this point on, samples from each attack were processed independently. For each attack, numerical features in the training and validation sets were first transformed using a logarithmic function to mitigate heavy-tailed distributions and were then standardized using the corresponding training-set mean and standard deviation. Discrete features were encoded using label encoding and treated as discrete variables throughout the modeling process.

## 4.3 Hyperparameter Optimization

To ensure a fair and reproducible comparison, hyperparameters for all DL-based generative models were selected using Bayesian optimization in MATLAB<sup>®</sup>, with a fixed budget of 50 objective function evaluations for each model. For all models, performance was monitored on the validation set, and the objective function corresponded to the standard formulation used for each model. Hyperparameters were optimized independently for each attack scenario, resulting in attack-specific model configurations reported in Table 2.

## 4.4 Baseline Generative Models

To assess the effectiveness of the LDM, we compared it against several data augmentation and generative modeling approaches, specifically SMOTE, VAE, GAN, and DM. Baseline models were trained exclusively on attack samples and were configured to handle mixed continuous and discrete features, with loss functions defined as a combination of continuous and discrete loss components.

Table 3: Configuration of baseline generative models.

Model	Hyperparameter	Value
SMOTE	Number of neighbors ( $k$ )	5
VAE	Encoder hidden widths	[128 64]
	Decoder hidden widths	[64 128]
	Activation function	ReLU
	Mini-batch size	128
	Max. epochs	300
	Optimizer	Adam
GAN	Generator hidden widths	[128 128]
	Discriminator hidden widths	[128 128]
	Activation function	ReLU
	Mini-batch size	128
	Max. epochs	300
	Optimizer	Adam
DM	Diffusion schedule	linear
	$[\beta_{start}, \beta_{end}]$	$[10^{-4}, 0.02]$
	Hidden width	256
	Time embedding dimension	128
	Mini-batch size	128
	Max. epochs	400
	Gradient clipping	1
	Optimizer	Adam

Table 4: Row-normalized CMs (%) for the different classifiers across attack types, with the best (bold) and second-best (italics) CMs highlighted. In each CM, rows correspond to true classes (Benign, Attack) and columns to predicted classes (Benign, Attack).

Attack	Baseline	SMOTE	VAE	GAN	DM	LDM
DDoS	$\begin{bmatrix} 100.0 & 0.0 \\ 69.9 & 30.1 \end{bmatrix}$	$\begin{bmatrix} 100.0 & 0.0 \\ 19.4 & 80.6 \end{bmatrix}$	$\begin{bmatrix} 100.0 & 0.0 \\ 7.5 & 92.5 \end{bmatrix}$	$\begin{bmatrix} 100.0 & 0.0 \\ 11.8 & 88.2 \end{bmatrix}$	$\begin{bmatrix} 100.0 & 0.0 \\ 6.5 & 93.5 \end{bmatrix}$	$\begin{bmatrix} \mathbf{100.0} & \mathbf{0.0} \\ \mathbf{1.1} & \mathbf{98.9} \end{bmatrix}$
Mirai	$\begin{bmatrix} 100.0 & 0.0 \\ 90.5 & 9.5 \end{bmatrix}$	$\begin{bmatrix} 100.0 & 0.0 \\ 51.0 & 49.0 \end{bmatrix}$	$\begin{bmatrix} 100.0 & 0.0 \\ 25.2 & 74.8 \end{bmatrix}$	$\begin{bmatrix} 100.0 & 0.0 \\ 28.6 & 71.4 \end{bmatrix}$	$\begin{bmatrix} 100.0 & 0.0 \\ 14.3 & 85.7 \end{bmatrix}$	$\begin{bmatrix} \mathbf{99.6} & \mathbf{0.4} \\ \mathbf{0.0} & \mathbf{100.0} \end{bmatrix}$
MitM	$\begin{bmatrix} 100.0 & 0.0 \\ 100.0 & 0.0 \end{bmatrix}$	$\begin{bmatrix} 99.7 & 0.3 \\ 78.3 & 21.7 \end{bmatrix}$	$\begin{bmatrix} 39.3 & 60.7 \\ 10.0 & 90.0 \end{bmatrix}$	$\begin{bmatrix} 35.5 & 64.5 \\ 8.3 & 91.7 \end{bmatrix}$	$\begin{bmatrix} \mathbf{81.5} & \mathbf{18.5} \\ \mathbf{26.7} & \mathbf{73.3} \end{bmatrix}$	$\begin{bmatrix} 75.7 & 24.3 \\ 31.7 & 68.3 \end{bmatrix}$

As a non-DL baseline, a SMOTE-based single-class data augmentation method was employed. Synthetic samples were generated by interpolating between randomly selected attack samples and their  $k$  nearest neighbors, computed using a mixed distance metric that combines Euclidean distance for continuous variables and Hamming distance for discrete variables. Continuous features were synthesized via linear interpolation, while discrete features were generated using a Bernoulli or categorical sampling strategy. A VAE was used as a latent-variable generative baseline. The VAE learned a probabilistic latent representation of attack traffic and generated synthetic samples by decoding latent variables drawn from a standard normal prior. A GAN was employed as a baseline generative model using an adversarial training framework. The generator mapped latent noise vectors to synthetic attack samples, while a discriminator was trained to distinguish real from generated data. Continuous features were generated directly by the generator, whereas discrete features were produced via sigmoid activations followed by Bernoulli or categorical sampling. A DM operating directly in data space was also used as a diffusion-based baseline. The model learned the distribution of attack samples by progressively corrupting real data with Gaussian noise over a predefined number of diffusion steps and training a neural network to reverse this process by predicting the added noise. The specific configurations of baseline models are reported in Table 3.

Table 5: F1-scores of the classifiers for each attack type, with the best (bold) and second-best (italics) results highlighted.

Attack	Baseline	SMOTE	VAE	GAN	DM	LDM
DDoS	0.46	0.89	0.96	0.94	<i>0.97</i>	<b>0.99</b>
Mirai	0.17	0.66	0.86	0.83	<i>0.92</i>	<b>0.99</b>
MitM	0.00	0.35	0.34	0.33	<b>0.53</b>	<i>0.44</i>

#### 4.5 Latent Diffusion Model Configuration Details

This subsection details the specific architectural choices, training settings, and implementation details adopted for the proposed LDM in the experimental evaluation. The deterministic AE used in the proposed LDM was implemented as a MLP with separate processing streams for numerical and categorical features. The encoder was composed of two parallel branches, i.e., a numerical branch with fully connected layers of widths 64 and 32, and a categorical branch with layers of widths 32 and 16, both using ReLU activations. The outputs of these branches were concatenated and passed through a shared fully connected layer with 32 units, followed by a linear projection to a latent space. The decoder mirrored this structure, consisting of a shared trunk with two fully connected layers of widths 32 and 64, followed by two output heads. The numerical head produced real-valued reconstructions of continuous features, while the categorical head outputted either Bernoulli or multinomial probabilities via a sigmoid or softmax activation, respectively. The AE was trained using a reconstruction loss composed of a MSE term for numerical variables and a categorical cross-entropy term for categorical variables. Training was performed using the Adam optimizer with batch size 128 for 200 epochs. After training the AE, all attack samples in the training set were encoded into latent representations. To ensure stable diffusion training, each latent dimension was normalized using the mean and standard deviation computed on the training set. These statistics were reused during sampling to denormalize generated latent vectors prior to decoding.

The DM operates exclusively in the normalized latent space. A linear variance schedule was employed with  $\beta_{\text{start}}$  set to  $10^{-4}$  and  $\beta_{\text{end}}$  set to 0.02. The reverse denoising network was implemented as a time-conditional MLP. A sinusoidal time embedding of dimension 128 was processed by a two-layer MLP and concatenated with the noisy latent input. The resulting representation was passed through a backbone network of optimized depth and hidden width 256, followed by a linear output layer predicting the clean latent representation. The DM was trained using a MSE loss between the predicted and true latent vectors. Optimization was performed using the Adam optimizer with gradient clipping, and training was carried out for 400 epochs. Synthetic samples were then generated by initializing the reverse diffusion process from Gaussian noise in latent space and iteratively applying the learned denoising transitions. The resulting latent samples were denormalized and decoded through the AE decoder to obtain synthetic samples. All generated samples corresponded exclusively to the attack class used during training.

## 5 Results

This section presents a comprehensive experimental evaluation of the proposed LDM for synthetic IoT attack data generation and its comparison with state-of-the-art approaches. The evaluation is structured to assess both the practical impact of data augmentation on IDS performance and the quality, fidelity, and efficiency of the generated samples.

First, a binary IDS (benign vs. attack) was constructed to assess whether augmenting minority-class attack data improves classification performance. The data available for evaluation, exhibiting attack-to-benign traffic ratios of approximately 1:5.30, 1:5.57, and 1:5.63 for the DDoS, Mirai, and MitM attack types, respectively, was split into training (80%) and test (20%) subsets. As the detection model, an ensemble classifier trained using the subspace method with 100 learning cycles was employed. To evaluate the impact of synthetic data augmentation, augmented training sets were created by adding synthetic attack samples generated by each generative model to the baseline training data until class balance was achieved. For each augmentation strategy, the classifier was retrained from scratch using the same configuration. All models were evaluated on the same imbalanced test set, and performance is summarized in Tables 4 and 5 in terms of row-normalized Confusion Matrices (CMs) and F1-scores, respectively.

Table 4 shows that training the classifier on the original, imbalanced datasets leads to a severe degradation in attack detection performance across all attack types. In the baseline setting, the classifier is strongly biased toward the majority class, yielding perfect benign traffic recognition but very low recall for attacks. This behavior is consistent with the extreme class imbalance and highlights the inability of the classifier to learn relevant attack patterns when minority samples are scarce. When the training set is balanced through the addition of synthetic attack samples, attack detection improves significantly for all augmentation strategies, as reflected by both the CMs in Table 4 and the corresponding F1-scores reported in Table 5. Among the evaluated augmentation methods, SMOTE consistently achieves the smallest performance gains across all attack types. Although it improves recall compared to the baseline, its performance remains clearly below that of DL-based models. This may be due to the fact that SMOTE generates new samples via linear interpolation between neighboring points, which tends to produce synthetic instances that add limited diversity to the attack class. In contrast, DL-based models provide synthetic samples that lead to higher detection



Table 6: Quantitative evaluation of synthetic samples generated for each attack type.

Metric	DDoS					Mirai					MitM				
	SMOTE	VAE	GAN	DM	LDM	SMOTE	VAE	GAN	DM	LDM	SMOTE	VAE	GAN	DM	LDM
MMD	<b>0.00</b>	0.10	0.01	0.05	0.06	<b>0.00</b>	0.16	0.01	0.08	0.08	<b>0.00</b>	0.51	0.01	0.02	0.02
Mean KL	<b>0.03</b>	0.65	0.53	0.78	0.70	<b>0.02</b>	1.13	0.48	0.55	0.61	<b>0.04</b>	1.18	0.55	0.62	0.67
MI error	<b>1.19</b>	7.59	3.86	4.66	4.67	<b>0.95</b>	10.50	4.91	5.24	5.01	<b>1.39</b>	9.16	5.16	4.17	4.05
Precision	0.99	0.84	<b>1.00</b>	<b>1.00</b>	0.99	<b>1.00</b>	0.87	0.89	0.82	0.99	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>
Recall	<b>1.00</b>	0.84	<b>1.00</b>	0.98	0.99	<b>1.00</b>	0.48	0.99	0.69	0.75	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>
$N_{\text{med}}$	0.47	<b>2.85</b>	0.62	0.83	0.88	0.28	<b>3.78</b>	0.82	0.83	1.96	1.49	<b>3.23</b>	1.98	1.66	2.33

performance. LDM achieves the best results for DDoS and Mirai attacks, reaching F1-scores of 0.99 in both cases. This indicates that latent diffusion is highly effective at enriching the attack sample space with realistic yet diverse instances that enhance classifier generalization. The DM operating directly in data space closely follows LDM performance for these attacks and yields the best results for the most challenging scenario, MitM, with an F1-score of 0.53, representing a notable improvement given the complete failure of the baseline classifier. LDM also performs competitively for MitM, achieving an F1-score of 0.44. VAE and GAN outperform SMOTE but remain inferior to diffusion-based approaches. While both methods increase the F1-score, their synthetic samples appear less effective at capturing the complex attack patterns, particularly for MitM traffic. These results demonstrate that the proposed LDM is well suited for synthetic IoT attack data generation, as it not only outperforms classical oversampling techniques but also achieves competitive or superior performance compared to state-of-the-art DL-based generative models.

To quantitatively assess the quality of the synthetic attack samples generated by LDM and to compare it with other augmentation techniques, we evaluated distributional similarity, sample diversity, and potential memorization effects. Specifically, Maximum Mean Discrepancy (MMD), mean marginal Kullback-Leibler (KL) divergence, Mutual Information (MI) error, generative precision and recall, and the median nearest-neighbor distance ( $N_{\text{med}}$ ) were computed between real attack samples and synthetic data generated by each model. These metrics were evaluated at three different augmentation ratios, specifically 50%, 100%, and 200% relative to the number of available real attack samples. Results are reported in Table 6 for the 100% augmentation ratio, as all ratios yielded similar metric values, suggesting that these quantitative metrics are mainly influenced by the quality of the generated samples rather than by the number of synthetic instances produced.

Table 6 shows that no single technique dominates across all metrics and attacks, highlighting the trade-offs that exist between distributional similarity, diversity, and memorization. SMOTE consistently achieves the lowest MMD and mean marginal KL values, indicating an almost perfect match between the empirical distributions of real and synthetic samples. It also yields the lowest MI error, suggesting strong preservation of feature dependencies. However, SMOTE also exhibits the lowest  $N_{\text{med}}$  across all attack types, which points to a considerable risk of memorization. This behavior is expected, as SMOTE generates new samples through linear interpolation between neighboring real instances, producing points that remain very close to the original data manifold and contribute limited novel information. GAN augmentation exhibits a similar, though less pronounced, trend. GAN achieves relatively low MMD and KL values, indicating good distributional alignment, however, its  $N_{\text{med}}$  values remain small, especially for DDoS and Mirai, suggesting that the generated samples are still concentrated near the real data and may not expand the attack space. This behavior is consistent with mode concentration effects commonly observed in GANs when modeling complex or fragmented distributions. In contrast, VAE shows the largest MMD and mean KL values, indicating a poorer match between real and synthetic distributions, and it also exhibits the highest MI error, suggesting that feature dependencies are less well preserved. However, VAE achieves the largest  $N_{\text{med}}$  values across all attack types, indicating minimal memorization and a strong tendency to generate novel samples. This reflects the smoothing effect of VAEs, where samples are drawn from a continuous latent distribution that may sacrifice fine-grained structure in favor of diversity. Diffusion-based models show a more favorable balance between these contrasting behaviors. Both methods achieve low MMD and KL values while maintaining moderate MI errors, indicating that they preserve both marginal distributions and feature dependencies reasonably well. With respect to  $N_{\text{med}}$ , LDM consistently obtains higher values than DM across all attack types, indicating that latent-space diffusion further enhances diversity while maintaining distributional fidelity. Generative precision and recall are generally high for all methods, confirming that most approaches generate samples that lie within the support of the real data distribution and provide adequate coverage. Results indicate that LDM offers a good compromise between realism, diversity, and robustness. While classical oversampling methods excel in distribution matching, they suffer from memorization, and VAEs favor novelty at the expense of fidelity. However, LDM effectively balances these competing objectives.

To qualitatively assess how synthetic samples affect the geometric structure of real attacks, Uniform Manifold Approximation and Projection (UMAP) projections of real and synthetic samples were analyzed in a three-dimensional embedding space. For each attack type and augmentation method, a combined dataset containing real attack samples and their corresponding synthetic counterparts was created, and a UMAP model was learned to visualize their joint manifold structure. Representative embeddings

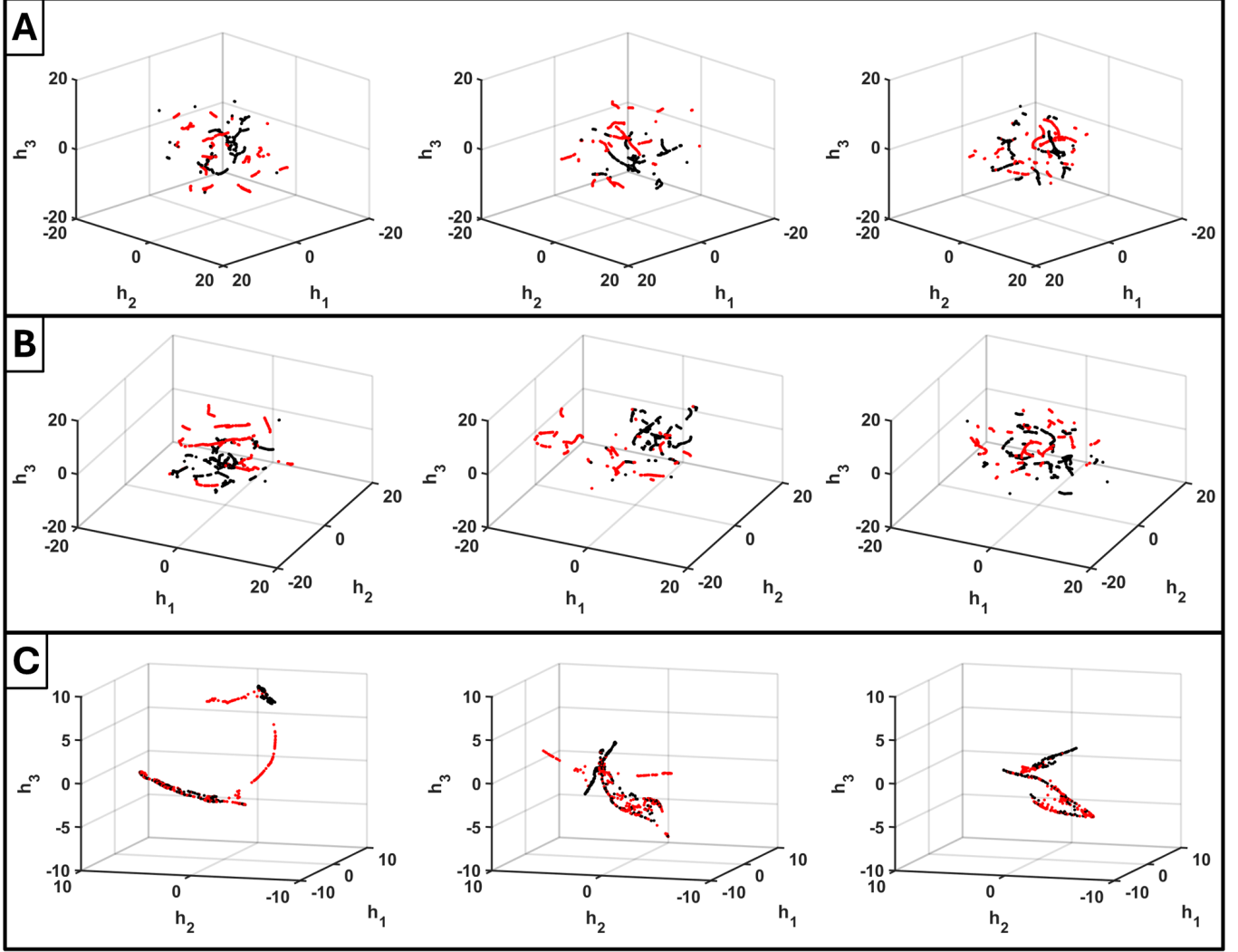


Figure 1: UMAP projections of real (black) and synthetic (red) samples generated by GAN (first column), DM (second column), and LDM (third column) for DDoS (panel A), Mirai (panel B), and MitM (panel C) attacks.

are shown in Fig. 1 for the DDoS (panel A), Mirai (panel B), and MitM (panel C) attack types, focusing on the manifolds obtained using GAN (first column), DM (second column), and LDM (third column). As expected, SMOTE-based manifolds exhibit synthetic samples that lie close to the real data, often nearly overlapping with them. This behavior reflects the interpolation-based nature of SMOTE and aligns with the low  $N_{\text{med}}$  values observed in Table 6, indicating limited diversity and a high degree of memorization. On the other hand, VAE-generated manifolds display similar behavior to GAN-generated manifolds, shown in Fig. 1.

For the DDoS attack type (Fig. 1, panel A), GAN- and DM-generated samples appear sparser than those produced by LDM. While GAN and DM introduce additional variability around the real data, several synthetic samples lie at the boundary of the real manifold, indicating partial coverage and, in the case of DM, the presence of off-manifold samples, particularly in the upper and left regions of the embedding. In contrast, LDM-generated samples effectively fill gaps in the manifold without collapsing onto existing points, suggesting that latent diffusion captures the local geometry of the attack space more accurately. For the Mirai attack (Fig. 1, panel B), off-manifold behavior becomes more evident for GAN and DM. These samples add diversity, but may change the intrinsic structure of the real data distribution. In contrast, LDM maintains synthetic samples within the real manifold boundaries, while still avoiding overlap with real points, indicating a favorable balance between fidelity and diversity. The MitM attack (Fig. 1, panel C) represents the most challenging scenario, as its real data manifold is fragmented. In this case, GAN-generated samples attempt to bridge disconnected regions of the manifold, which alters the original geometric structure. DM exhibits a similar tendency, with several synthetic samples deviating from the real clusters. In contrast, LDM-generated samples are distributed across the fragmented regions of the real manifold, reinforcing existing structures.

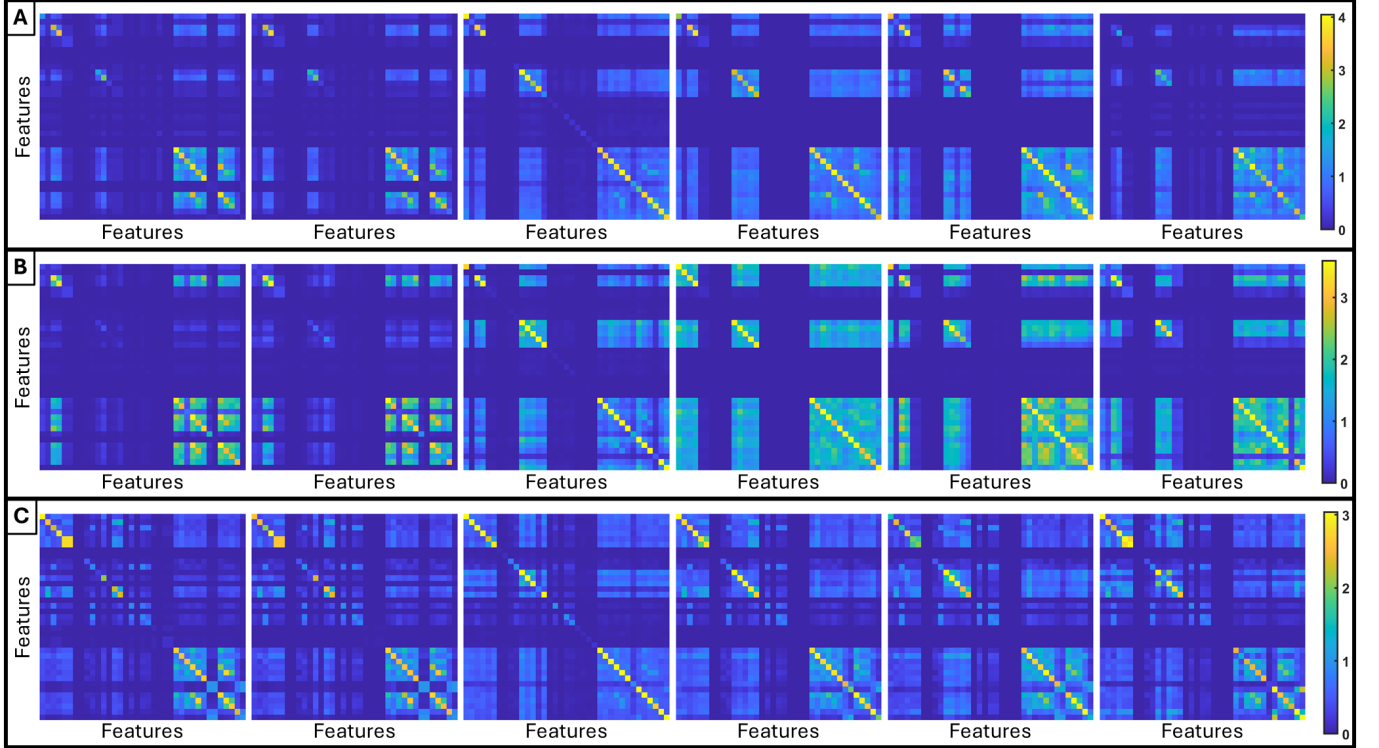


Figure 2: MI matrices for DDoS (panel A), Mirai (panel B), and MitM (panel C) attack features. Columns show the MI matrices computed from the full real attack set and from a reduced subset of 20 real attack samples augmented with synthetic data generated using SMOTE, VAE, GAN, DM, and LDM, from left to right. Color intensity reflects the strength of pairwise feature dependencies.

To further analyze whether the statistical dependencies among features are preserved by the different augmentation methods, MI matrices were computed for each attack type and generative method. For each attack (panels A, B and C in Fig. 2 correspond to DDoS, Mirai, and MitM attacks, respectively), the MI matrix of the full real attack dataset was first computed (first column of Fig. 2). Then, a reduced subset of 20 randomly selected real attack samples was augmented with synthetic samples generated using SMOTE, VAE, GAN, DM, and LDM until matching the size of the original attack dataset. The MI matrices computed from these augmented datasets are shown in columns 2 to 6 on Fig. 2.

As expected, augmenting the reduced real attack subsets with SMOTE yields MI matrices that almost perfectly match those computed from the full real attack datasets for all attack types. This behavior is consistent with the very low MI error values reported in Table 6. Since SMOTE generates synthetic samples through linear interpolation between existing instances, it largely preserves the original feature relationships but introduces limited diversity. In contrast, deep generative models introduce greater variability, which is reflected in larger differences in the MI matrices. For all attack types, GAN- and DM-augmented datasets exhibit more noticeable differences with respect to the real MI matrices, particularly in the attenuation or distortion of localized dependency patterns. Between these two methods, DM generally preserves feature dependencies slightly better than GAN, but still shows regions where correlations are weakened or displaced. VAE-based augmentation produces MI matrices that more closely resemble the real attack dependencies than GAN and DM in several regions. In particular, finer dependency structures, visible in the central feature blocks and in the upper-left regions of the matrices, are partially retained, whereas they are absent or blurred for GAN and DM. Among all approaches, LDM consistently yields MI matrices that best balance fidelity and diversity. Across all attack types, LDM preserves the global dependency structure while maintaining finer-grained details that are lost or distorted by the other DL methods. This qualitative evidence supports the quantitative results and highlights the ability of latent diffusion to generate diverse synthetic samples while preserving meaningful feature dependencies.

Finally, the computational efficiency of the different generative methods was evaluated by measuring the time required to generate synthetic samples. Since data augmentation is typically performed offline but may involve large numbers of samples, sampling efficiency remains an important practical consideration. Generation times were measured for  $10^2$ ,  $10^3$ , and  $10^4$  synthetic samples, with each experiment repeated five times. The mean and standard deviation of the sampling time are reported in Table 7. As expected, GAN is the fastest across all sample sizes, requiring only a few milliseconds even for generating  $10^4$  samples. This reflects the fact that, once trained, GAN sampling consists of a single forward pass through the generator network. VAE also

Table 7: Time required to generate synthetic samples (mean  $\pm$  std) by each method for increasing sample sizes, with the fastest (bold) and second-fastest (italics) methods highlighted.

	10 <sup>2</sup> samples	10 <sup>3</sup> samples	10 <sup>4</sup> samples
SMOTE	137.17 $\pm$ 3.12 ms	139.45 $\pm$ 5.49 ms	154.07 $\pm$ 2.86 ms
VAE	<i>5.97 <math>\pm</math> 1.56 ms</i>	<i>9.19 <math>\pm</math> 1.97 ms</i>	<i>54.85 <math>\pm</math> 9.44 ms</i>
GAN	<b>2.62 <math>\pm</math> 0.67 ms</b>	<b>4.00 <math>\pm</math> 0.78 ms</b>	<b>4.85 <math>\pm</math> 0.87 ms</b>
DM	487.16 $\pm$ 6.59 ms	496.01 $\pm$ 3.08 ms	553.86 $\pm$ 9.31 ms
LDM	213.24 $\pm$ 4.45 ms	259.06 $\pm$ 7.81 ms	411.23 $\pm$ 8.34 ms

exhibits low sampling times, although its cost increases more noticeably with the number of generated samples due to stochastic decoding and sampling operations. These results confirm that latent-variable models relying on direct decoding are highly efficient at inference time. SMOTE shows higher and comparatively stable generation times across all sample sizes. This behavior arises because its computational cost is dominated by nearest-neighbor searches in the original dataset rather than by the number of synthetic samples generated. Diffusion-based models are the most computationally demanding. DM exhibits the highest sampling times, reflecting the need to perform a large number of iterative denoising steps in the original high-dimensional feature space. In contrast, LDM substantially reduces sampling cost compared to DM for all sample sizes. For 10<sup>4</sup> synthetic samples, LDM achieves an approximate 25% reduction in generation time, highlighting the efficiency gains obtained by performing diffusion in a lower-dimensional space while retaining the benefits of diffusion-based generation.

## 6 Discussion

Our results confirm that the extreme class imbalance inherent to IoT traffic datasets significantly limits the ability of supervised IDSs to learn meaningful attack patterns. When trained on the original imbalanced data, the classifier was dominated by the overwhelming presence of benign traffic and tended to converge toward a decision rule that favors benign predictions, resulting in high benign recognition rates but limited attack detection capability. This behavior is a well-known consequence of minority-class under-representation, where the scarcity of attack samples prevents the classifier from learning well-defined discriminative decision boundaries. By balancing the training set through synthetic attack data augmentation, the learning dynamics were altered. Attack samples became sufficiently represented to meaningfully influence the training process, allowing decision boundaries to move away from the benign majority and better capture the structure of the attack class. As a result, attack detection performance improved consistently across all evaluated attack types. However, the magnitude of this improvement was not uniform and depended on the intrinsic complexity of the attack. Attack types that are more different from benign traffic, such as DDoS and Mirai, which exhibit pronounced flooding behavior, benefited the most from data augmentation. In contrast, attacks whose traffic patterns are more similar to benign activity, such as MitM, remained more challenging to detect even after balancing. This observation highlights that data augmentation substantially mitigates, but does not completely eliminate, the difficulty posed by subtle and complex attack behaviors.

Beyond the general benefit of balancing the training set, the results reveal clear differences among augmentation strategies in how effectively they enrich the minority attack class. Classical oversampling with SMOTE consistently preserved marginal data statistics and pairwise feature dependencies, as reflected by its very low MMD values, mean KL divergence, and MI error. However, because SMOTE relies on linear interpolation between neighboring samples, it introduced limited diversity and exhibited a memorization-like behavior, as indicated by the smallest  $N_{\text{med}}$  values. Consequently, despite its excellent distributional similarity metrics, SMOTE provided only modest gains in IDS performance, as it largely replicated existing attack patterns rather than expanding the effective decision space available to the classifier. At the opposite end, VAE generated highly diverse samples, evidenced by large  $N_{\text{med}}$  values, but struggled to preserve the true data distribution and feature dependencies, particularly for more complex attacks such as MitM. This behavior is consistent with the high MMD, KL divergence, and MI error values observed, and can be due to the Gaussian latent prior and the tendency of VAEs to oversmooth multimodal or fragmented data manifolds. GAN offered an intermediate behavior, improving sample realism and achieving lower MMD and KL divergence than VAE. However, GANs adversarial training dynamics can lead to irregular coverage of the data manifold and occasional mode collapse, which manifested as off-manifold samples in low-dimensional projections and distortions in MI structure. The DM operating directly in data space addressed several of these limitations by explicitly modeling the data distribution through iterative denoising, resulting in high-fidelity samples and strong gains in classification performance. However, this came at the cost of increased computational complexity, as the large number of diffusion steps required in the high-dimensional feature space led to slow sampling, and occasional off-manifold behavior is still observed. The proposed LDM effectively combined the strengths of diffusion-based generation with the efficiency and regularization provided by a learned latent representation. By performing diffusion in a compact and well-conditioned latent space, LDM not only achieved a favorable balance between realism, coverage, diversity, and dependency preservation, but also significantly reduced sampling time compared to data-space diffusion.

This improved efficiency, together with the quality of the generated samples, is consistently reflected across quantitative metrics, manifold visualizations, MI analyses, and downstream IDS performance.

DDoS and Mirai attacks exhibit comparatively structured and homogeneous traffic patterns, which makes them easier to model and distinguish from benign traffic. DDoS flows are typically characterized by consistent packet-level statistics, such as high request rates, repetitive communication patterns, and limited protocol variability, resulting in dense and well-connected manifolds in feature space. Similarly, Mirai attacks, while more diverse than DDoS, are driven by automated botnet behavior that produces recurring scanning and exploitation signatures, leading to attack manifolds that remain relatively compact and continuous. These properties facilitate both generative modeling and downstream classification. In contrast, MitM attacks represent a substantially more challenging scenario. MitM traffic is inherently sparse, heterogeneous, and highly context-dependent, often closely resembling benign communication while introducing only slight changes in packet timing, routing behavior, or protocol semantics. Therefore, the associated feature distributions are fragmented and discontinuous, forming few isolated regions in latent space rather than a single cluster. Diffusion-based models are better suited to capturing such complex and multimodal distributions, yet they also exhibited limitations when modeling MitM traffic. DM achieved the highest F1-score for MitM detection, outperforming the proposed LDM, suggesting that direct diffusion in feature space can be advantageous for improving classification performance in this difficult setting. However, UMAP projections revealed that data-space diffusion occasionally produces samples that lie outside the real attack manifold, indicating potential geometric inconsistencies with the true data distribution. In contrast, LDM generated synthetic samples that more coherently reconstructed the fragmented structure of the MitM attack manifold, distributing samples across disconnected submanifolds while remaining largely aligned with the geometry of the real data. This highlights a trade-off between maximizing detection performance and preserving the intrinsic structure of complex attack distributions.

Despite the promising results obtained in this study, several limitations should be acknowledged. First, all experiments were conducted on a single IoT attack dataset and focused on three representative attack types. Although these attacks exhibit distinct statistical and geometric characteristics, the generalizability of the proposed approach to other datasets, network environments, and attack scenarios remains to be systematically evaluated. Second, the experimental analysis was restricted to binary classification settings. While this choice is appropriate for isolating and assessing the impact of minority-class data augmentation, it does not fully capture the complexity of real-world IDSs, where multiple attack classes may coexist and partially overlap. Future work will therefore focus on extending the proposed framework to cross-dataset evaluation scenarios in order to assess robustness under domain shifts and varying traffic distributions. In addition, integrating latent diffusion-based data augmentation into multi-class intrusion detection pipelines represents a natural and practically relevant extension of this work. Finally, further improvements in sampling efficiency, as well as the exploration of online or incremental augmentation strategies, may enhance the applicability of LDMs in large-scale and real-time IoT security deployments.

## 7 Conclusions

The severe class imbalance that characterizes real-world IoT intrusion detection datasets poses a fundamental challenge for data-driven security models, as minority attack patterns are often underrepresented and poorly learned by conventional classifiers. Existing data augmentation strategies either rely on simple interpolation mechanisms or struggle to capture the complex and heterogeneous distributions of network traffic, thus limiting their effectiveness in enriching underrepresented attack classes. In this work, we address this challenge by proposing a latent diffusion-based framework for synthetic IoT attack data generation. By combining a deterministic AE with a denoising diffusion process operating in a compact latent space, the proposed LDM overcomes the scalability limitations of data-space diffusion while preserving the ability to model complex, multimodal attack distributions. This design enables the generation of diverse and realistic synthetic samples from limited real attack data, without requiring access to benign traffic during generation. Experimental results demonstrate that latent diffusion consistently improves attack detection performance when used for data augmentation. Compared to classical oversampling techniques and established deep generative models, the proposed approach achieves a favorable balance between distributional similarity, sample diversity, dependency preservation, and computational efficiency. Specifically, latent diffusion yields substantial gains in IDS performance for DDoS and Mirai attacks and provides meaningful improvements for the more challenging MitM scenario, as confirmed by quantitative metrics, manifold analyses, MI preservation, and classification results. Therefore, this study highlights LDMs as an effective and scalable solution for synthetic tabular data generation in IoT security. Beyond intrusion detection, the proposed framework offers a generalizable approach for addressing data scarcity and class imbalance in other cybersecurity and network analytics applications, paving the way for more robust and data-efficient learning systems.

## Acknowledgments

This work was partially supported by CyberFold (ETD202300129) from the European Union (EU NextGeneration, PRTR) through INCIBE; by LATENTIA (PID2022-140786NB-C31) and HERMES (PID2023-152331OA-I00), funded by MICIU/AEI/10.13039/501100011033 and by FEDER/UE; and by the ELLIS Madrid Node (Autonomous Community of Madrid).

## References

- Yusuf Perwej, Kashiful Haq, Firoj Parwej, and Mumdouh M. Mohamed Hassan. The internet of things (IoT) and its application domains. *International Journal of Computer Applications*, 182(49):36–49, 2019.
- Ritika Lohiya and Ankit Thakkar. Application domains, evaluation data sets, and research challenges of IoT: A systematic review. *IEEE Internet of Things Journal*, 8(11):8774–8798, 2021.
- Muhammad Shafiq, Zhaoquan Gu, Omar Cheikhrouhou, Wajdi Alhakami, and Habib Hamam. The rise of “Internet of Things”: Review and open research issues related to detection and prevention of IoT-based security attacks. *Wireless Communications and Mobile Computing*, 2022(8669348):1–12, 2022.
- Ansam Khraisat, Iqbal Gondal, Peter Vamplew, and Joarder Kamruzzaman. Survey of intrusion detection systems: techniques, datasets and challenges. *Cybersecurity*, 2(20):1–22, 2019.
- Alex Shenfield, David Day, and Aladdin Ayeshe. Intelligent intrusion detection systems using artificial neural networks. *ICT Express*, 4(2):95–99, 2018.
- Francisco M. Melgarejo-Meseguer, Andrés Lorenzo-Bleda, Sergio Eduardo-Abbenante, Francisco Javier Gimeno-Blanes, Estrella Everss-Villalba, Sergio Muñoz-Romero, José Luis Rojo-Álvarez, and Rafael Maestre Ferriz. Anomaly detection from low-dimensional latent manifolds with home environmental sensors. *IEEE Internet of Things Journal*, 11(8):13411–13423, 2021.
- Haoqi Huang, Ping Wang, Jianhua Pei, Jiacheng Wang, Shahen Alexanian, and Dusit Niyato. Deep learning advancements in anomaly detection: A comprehensive survey. *IEEE Internet of Things Journal*, 12(21):44318–44342, 2025.
- Barjinder Kaur, Sajjad Dadkhah, Farzaneh Shoeleh, Euclides Carlos Pinto-Neto, Pulei Xiong, Shahrear Iqbal, Philippe Lamontagne, Suprio Ray, and Ali A. Ghorbani. Internet of things (IoT) security dataset evolution: Challenges and future directions. *Internet of Things*, 22(100780):1–23, 2023.
- Christin Alex, Giselle Creado, Wesam Almobaideen, Orieb Abu Alghanam, and Maha Saadeh. A comprehensive survey for IoT security datasets taxonomy, classification and machine learning mechanisms. *Computers & Security*, 132(103283):1–21, 2023.
- Elias Dritsas and Maria Trigka. A survey on cybersecurity in IoT. *Future Internet*, 17(30):1–32, 2025.
- Huan Ma, Wan Zhang, Dalong Zhang, and Baozhan Chen. An IoT intrusion detection framework based on feature selection and large language models fine-tuning. *Scientific Reports*, 15(21158):1–21, 2025.
- S Kumar Reddy Mallidi and Rajeswara Rao Ramisetty. Optimizing intrusion detection for IoT: A systematic review of machine learning and deep learning approaches with feature selection and data balancing. *WIREs Data Mining and Knowledge Discovery*, 15(2):1–41, 2025.
- Md Mahbubur Rahman, Shaharia Al Shakil, and Mizanur Rahman Mustakim. A survey on intrusion detection system in IoT networks. *Cyber Security and Applications*, 3(100082):1–19, 2025.
- Nitesh V. Chawla, Kevin W. Bowyer, Lawrence O. Hall, and W. Philip Kegelmeyer. SMOTE: synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16(1):321–357, 2002.
- Nur Athirah Azhar, Muhammad Syafiq Mohd Pozi, Aniza Mohamed Din, and Adam Jatowt. An investigation of SMOTE based methods for imbalanced datasets with data complexity analysis. *IEEE Transactions on Knowledge and Data Engineering*, 35(7):6651–6672, 2022.
- Shuo Zhang, Jinyi Chen, Jiayuan Chen, Xiaofei Chen, and Hejiao Huang. Data imputation in IoT using spatio-temporal variational auto-encoder. *Neurocomputing*, 529:23–32, 2023.
- Kunfeng Wang, Chao Gou, Yanjie Duan, Yilun Lin, Xinhua Zheng, and Fei-Yue Wang. Generative adversarial networks: Introduction and outlook. *IEEE/CAA Journal of Automatica Sinica*, 4(4):588–598, 2017.
- Alex X. Wang, Stefanka S. Chukova, Colin R. Simpson, and Binh P. Nguyen. Challenges and opportunities of generative models on tabular data. *Applied Soft Computing*, 166(112223):1–15, 2024.
- Xiang Li, Yang Peng, and Hao Zheng. Research and analysis of VAE, GAN, and diffusion generation models. *Science and Technology of Engineering, Chemistry and Environmental Protection*, 1(6):1–5, 2025a.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Proceedings of the 34th International Conference on Neural Information Processing Systems (NeurIPS)*, pages 6840–6851, 2020.
- Zhong Li, Qi Huang, Lincen Yang, Jiayang Shi, Zhao Yang, Niki van Stein, Thomas Back, and Matthijs van Leeuwen. Diffusion models for tabular data: Challenges, current progress, and future directions. *arXiv*, 2025b.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the 35th IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022.

- Timur Sattarov, Marco Schreyer, and Damian Borth. FinDiff: Diffusion models for financial tabular data generation. In *Proceedings of the 4th ACM International Conference on AI in Finance (ICAIF)*, pages 64–72, 2023.
- Mario Villaizán-Vallelado, Matteo Salvatori, Carlos Segura, and Ioannis Arapakis. Diffusion models for tabular data imputation and synthetic data generation. *ACM Transactions on Knowledge Discovery from Data*, 19(125):1–32, 2025.
- Chaoyi Zhu, Jiayi Tang, Juan F. Pérez, Marten van Dijk, and Lydia Y. Chen. DP-TLDM: Differentially private tabular latent diffusion model. In *Proceedings of the 20th International Conference on Availability, Reliability and Security (ARES)*, pages 337–357, 2025.
- Yan Naung Soe, Paulus Insap Santosa, and Rudy Hartanto. DDoS attack detection based on simple ANN with SMOTE for IoT environment. In *Proceedings of the 4th International Conference on Informatics and Computing (ICIC)*, pages 1–5, 2019.
- Pandit Byomakesha Dash, Janmenjoy Nayak, Bighnaraj Naik, Etuari Oram, and SK Hafizul Islam. Model based IoT security framework using multiclass adaptive boosting with SMOTE. *Security and Privacy*, 3(5):1–15, 2020.
- Lucas Pinheiro-Cinelli, Matheus Araújo-Marins, Eduardo Antônio Barros da Silva, and Sérgio Lima-Netto. Variational autoencoder. In *Variational Methods for Machine Learning with Applications to Deep Networks*, chapter 5, pages 111–149. Springer, Cham, Berlin, DE, 2021.
- Chang Liu, Ruslan Antypenko, Iryna Sushko, and Oksana Zakharchenko. Intrusion detection system after data augmentation schemes based on the VAE and CVAE. *IEEE Transactions on Reliability*, 71(2):1000–1100, 2022.
- Ying Zhang and Qiang Liu. On IoT intrusion detection based on data augmentation for enhancing learning on unbalanced samples. *Future Generation Computer Systems*, 133:213–227, 2022.
- Cheolhee Park, Jonghoon Lee, Youngsoo Kim, Jong-Geun Park, Hyunjin Kim, and Dowon Hong. An enhanced AI-based network intrusion detection system using generative adversarial networks. *IEEE Internet of Things Journal*, 10(3):2330–2345, 2023.
- Saifur Rahman, Shantanu Pal, Shubh Mittal, Tisha Chawla, and Chandan Karmakar. SYN-GAN: A robust intrusion detection system using GAN-based synthetic data for IoT security. *Internet of Things*, 26(101212):1–13, 2024.
- Chiara Camerota, Lorenzo Pappone, Tommaso Pecorella, and Flavio Esposito. Addressing data security in IoT: Minimum sample size and denoising diffusion models for improved malware detection. In *Proceedings of the 20th International Conference on Network and Service Management (CNSM)*, pages 1–7, 2024.
- Peng Wang, Yafei Song, and Xiaodan Wang. DIFT: A diffusion-transformer for intrusion detection of IoT with imbalanced learning. *Journal of Network and Systems Management*, 33(48):1–24, 2025.
- Akim Kotelnikov, Dmitry Baranchuk, Ivan Rubachev, and Artem Babenko. TabDDPM: Modelling tabular data with diffusion models. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, pages 17564–17579, 2023.
- Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, pages 2319–2329, 2021.
- Tero Karras, Miika Aittala, Samuli Laine, and Timo Aila. Elucidating the design space of diffusion-based generative models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems (NeurIPS)*, pages 26565–26577, 2022.
- Euclides Carlos Pinto-Neto, Sajjad Dadkhahand Raphael Ferreira, Alireza Zohourian, Rongxing Lu, and Ali A. Ghorbani. CICIOT2023: A real-time dataset and benchmark for large-scale attacks in IoT environment. *Sensors*, 23(13):1–26, 2023.
- Krushang Sonar and Hardik Upadhyay. A survey: DDOS attack on internet of things. *International Journal of Engineering Research and Development*, 10(11):58–63, 2014.
- Syed Ghazanfar Abbas, Fabiha Hashmat, Ghalib A. Shah, and Kashif Zafar. Generic signature development for IoT botnet families. *Forensic Science International: Digital Investigation*, 38(301224):1–14, 2021.
- Hamidreza Fereidouni, Olga Fadeitseva, and Mehdi Zalai. IoT and man-in-the-middle attacks. *Forensic Science International: Digital Investigation*, 8(2):1–19, 2025.