

# Synthetic Network Traffic Data Generation: A Comparative Study

Dure Adan Ammara  
Department of Computer Science  
Blekinge Institute of Technology  
Karlskrona, Sweden  
dure.adan.ammara@bth.se

Jianguo Ding  
Department of Computer Science  
Blekinge Institute of Technology  
Karlskrona, Sweden  
jianguo.ding@bth.se

Kurt Tutschku  
Department of Computer Science  
Blekinge Institute of Technology  
Karlskrona, Sweden  
kurt.tutschku@bth.se

**Abstract**—The generation of synthetic network traffic data is essential for network security testing, machine learning model training, and performance analysis. However, existing methods for synthetic data generation differ significantly in their ability to maintain statistical fidelity, utility for classification tasks, and class balance. This study presents a comparative analysis of twelve synthetic network traffic data generation methods, encompassing non-AI (statistical), classical AI, and generative AI techniques. Using NSL-KDD and CIC-IDS2017 datasets, we evaluate the fidelity, utility, class balance, and scalability of these methods under standardized performance metrics. Results demonstrate that GAN-based models, particularly CTGAN and CopulaGAN, achieve superior fidelity and utility, making them ideal for high-quality synthetic data generation. Statistical methods such as SMOTE and Cluster Centroid effectively maintain class balance but fail to capture complex traffic structures. Meanwhile, diffusion models exhibit computational inefficiencies, limiting their scalability. Our findings provide a structured benchmarking framework for selecting the most suitable synthetic data generation techniques for network traffic analysis and cybersecurity applications.

**Index Terms**—Synthetic tabular Data, Network traffic, NSL-KDD, CIC-IDS, Mutual information, Generative Adversarial Networks (GANs), Diffusion models, generative AI

## I. INTRODUCTION

The increasing reliance on data-driven decision-making in networking and security has intensified the need for high-quality network traffic data. This data plays a fundamental role in various applications, including network performance analysis, security research, and the development of machine learning (ML) models [1]. However, collecting and utilizing real-world network traffic data presents significant challenges. Issues such as privacy concerns, limited availability, and the high cost of manual labeling restrict access to large, diverse, and representative datasets [2].

To overcome these limitations, synthetic network traffic data generation has emerged as a promising alternative [3]–[7]. Synthetic datasets can simulate network traffic patterns while mitigating ethical, legal, and confidentiality concerns, enabling researchers and practitioners to develop, evaluate, and deploy data-driven solutions in a more scalable manner [8]. Despite this potential, existing methods for synthetic network traffic generation vary widely in their effectiveness, and no comprehensive comparison exists to guide the selection of

the most suitable approach for synthetic network traffic data generation.

### A. The Need for Synthetic Network Traffic Data

Synthetic network traffic data generation is important for several reasons:

- 1) **Privacy and Data Sensitivity:** Due to regulatory restrictions (e.g., GDPR) and commercial confidentiality, real network traffic data is often inaccessible [9]. Synthetic data provides a privacy-preserving alternative without exposing sensitive user information [10].
- 2) **Data Imbalance and Scarcity:** Real-world network traffic datasets are highly imbalanced, with benign traffic dominating over malicious activity, which can limit the effectiveness of ML-based models [11]. Synthetic data can address this by generating balanced datasets for better model training and evaluation.
- 3) **Scenario Diversity:** Real-world datasets may not cover a broad range of traffic conditions. Synthetic data generation techniques allow for the creation of diverse traffic scenarios, including rare or unseen network behaviors, which are crucial for testing robust ML models [12], [13].
- 4) **Scalability and Adaptability:** Unlike real data collection, which requires extensive monitoring infrastructure, synthetic data can be generated on demand to match specific analytical or experimental needs, reducing reliance on costly and time-consuming data collection processes [14].

### B. Challenges in Synthetic Network Traffic Generation

Despite its advantages, generating high-quality synthetic network traffic data remains challenging:

- 1) **Ensuring Realism:** Many synthetic data generation methods struggle to replicate the statistical and temporal properties of real network traffic, limiting their practical applicability [15].
- 2) **Variability and Generalization:** Some generative models suffer from mode collapse, producing limited traffic variations, which reduces their effectiveness in representing diverse network conditions [16]–[18].

- 3) **Computational Complexity:** State-of-the-art generative methods, such as Generative Adversarial Networks (GANs) and diffusion models, require significant computational resources, which can hinder their scalability for real-world deployment [19].
- 4) **Class Imbalance:** Many real-world network datasets exhibit severe class imbalances, where attack or rare traffic classes are underrepresented. Existing synthetic data generation techniques often fail to address this imbalance effectively, leading to biased models that struggle with anomaly detection and classification [20].

### C. Motivation

Generating synthetic network traffic data plays a critical role in network analysis, security testing, and developing machine learning models. While various methods exist for generating synthetic traffic, ranging from non-AI (statistical techniques) to AI (classical and generative) approaches, their effectiveness varies significantly [21], [22]. Prior studies have explored these methods individually, yet their comparative strengths and weaknesses remain unclear.

One of the key limitations in existing research is the lack of a systematic evaluation of synthetic data generation techniques across different methodological categories. Most studies focus on a single class of methods without directly comparing their ability to generate high-fidelity and utility-preserving data.

Another challenge lies in handling class imbalance in network traffic datasets. Many real-world network environments exhibit highly skewed distributions, where benign traffic significantly outnumbers attack or anomaly instances. While synthetic data generation has been proposed as a solution, current methods often fail to maintain class balance without compromising data fidelity [23].

To address these limitations, this study systematically evaluates and compares multiple synthetic network traffic generation methods across non-AI (statistical) and AI (classical and generative) approaches. By assessing their fidelity, utility, and class-balancing capabilities, this research provides valuable insights into selecting the most suitable synthetic data generation method for different network traffic applications.

### D. Objectives

This study aims to systematically evaluate synthetic network traffic data generation methods by addressing the following research questions:

- 1) **Fidelity and Realism:** How well do different synthetic data generation techniques preserve real network traffic's statistical properties and temporal dynamics?
- 2) **Utility for Machine Learning Tasks:** How effectively can synthetic data generated by different methods improve the performance of ML-based network traffic classification models?
- 3) **Class Balance and Data Diversity:** Which methods most effectively address class imbalance in network traffic datasets while maintaining statistical similarity?

- 4) **Scalability and Computational Efficiency:** What are the computational trade-offs associated with each method, and how do they impact scalability for large-scale simulations?

### E. Contributions

This study presents a systematic benchmarking of synthetic network traffic generation methods across statistical (non-AI) and AI (classical and generative) AI. Our key contributions are:

1) **Comprehensive Comparative Analysis:** We systematically evaluate twelve prominent synthetic data generation methods using standardized metrics, including fidelity, utility for machine learning models, class balance, and computational efficiency. This study provides a structured comparison under a unified experimental framework, addressing the lack of direct performance comparisons in prior work.

2) **Empirical Insights into Trade-offs Between Fidelity, Utility, and Computational Cost:** Our analysis reveals that generative AI models achieve superior fidelity but require higher computational resources and often struggle with class balance. In contrast, statistical methods effectively handle class imbalance but fail to preserve complex traffic patterns. These findings provide practical recommendations for selecting the most suitable method based on specific application needs.

3) **Evaluation of Class Imbalance Handling in Synthetic Network Traffic Generation:** We analyze how different methods address the issue of class imbalance, a persistent challenge in real-world network traffic datasets. Our results highlight the effectiveness of various techniques in generating balanced yet realistic traffic distributions, contributing to the development of more reliable synthetic datasets for machine learning applications.

The source code for this study is publicly available at [https://github.com/AdenRajput/Comparative\\_Analysis.git](https://github.com/AdenRajput/Comparative_Analysis.git).

### F. Paper Organization

The rest of this paper is structured as follows:

- **Section II** reviews existing synthetic data generation methods and highlights gaps in prior research.
- **Section III** describes the non-AI (statistical) and AI (classical and generative) based synthetic tabular data generation methods.
- **Section IV** introduces the evaluation metrics used to assess fidelity, utility, and class balance.
- **Section V** presents the experimental setup, including datasets, pre-processing steps, and model configurations.
- **Section VI** provides a detailed comparative analysis of results across different methods.
- **Section VII** outlines future research directions and potential improvements in synthetic network traffic generation.

## II. RELATED WORK

The generation of synthetic network traffic data has gained significant attention due to challenges associated with using

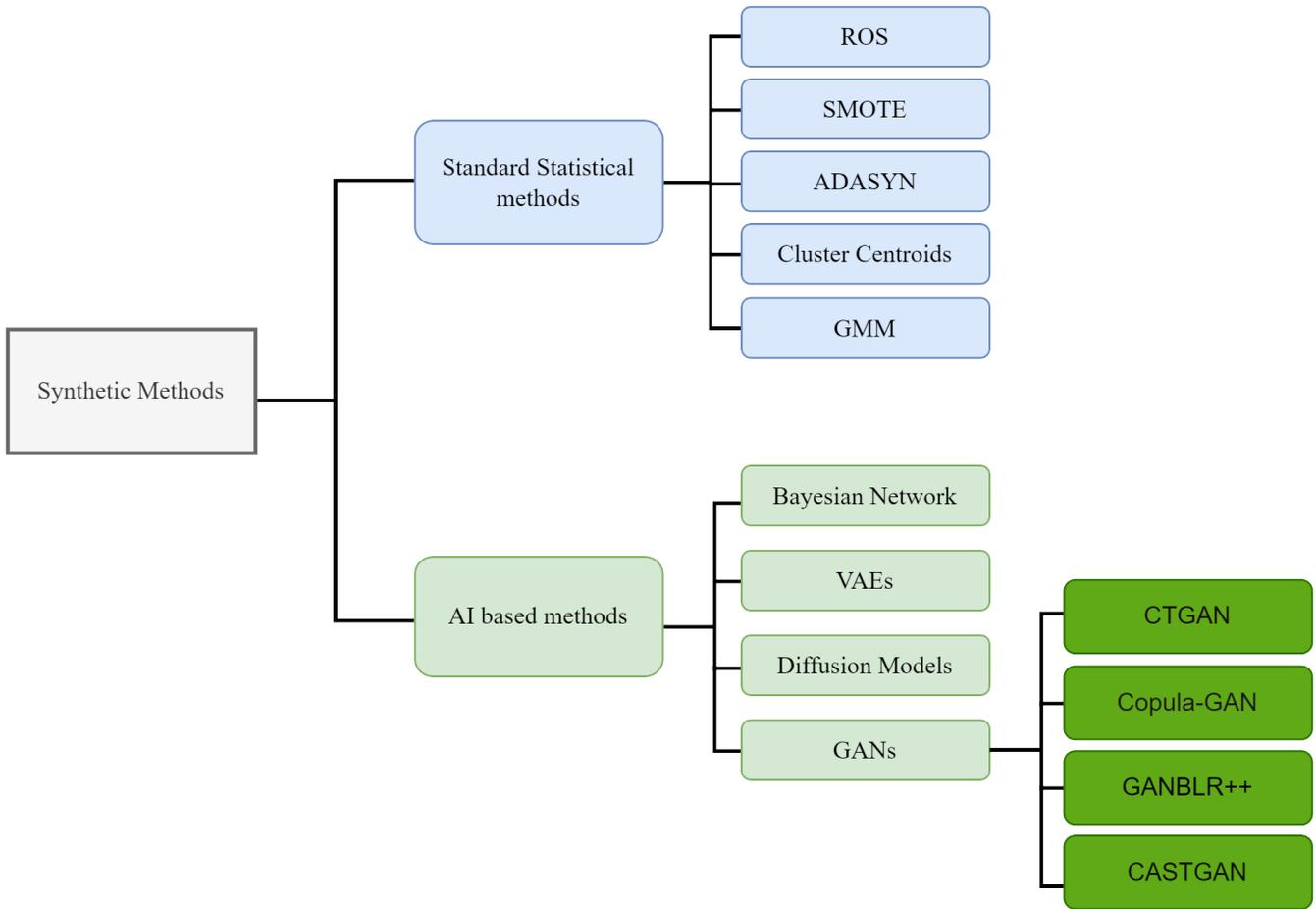


Fig. 1. Methods for generating synthetic tabular data

real-world datasets, such as privacy concerns, ethical limitations, and data imbalance. Existing approaches for synthetic data generation can broadly be classified into Non-AI methods and State-of-the-Art (SOTA) AI methods.

- **Non-AI Methods for Synthetic Data Generation:**

Non-AI approaches primarily focus on statistical techniques and rule-based models for generating synthetic tabular data. Common methods include ROS (Random Over-Sampling), SMOTE (Synthetic Minority Over-sampling Technique), ADASYN (Adaptive Synthetic Sampling), GMM (Gaussian Mixture Models), and Cluster Centroids (CC) [22]. These methods are often applied in scenarios addressing class imbalance rather than comprehensive traffic simulation [36]. For instance, SMOTE and ADASYN have been used to oversample minority classes within network traffic datasets and are frequently compared with modern AI methods such as Variational Autoencoders (VAEs) and GANs (Generative Adversarial Networks) [37].

Hybrid approaches such as SMOTE-GAN, COPULA-GAN, and Bayesian GAN (GANBLR) [38], [39] have

also been explored, combining statistical and generative AI techniques to enhance synthetic data fidelity [40]–[42]. Despite their effectiveness in managing class imbalance, these approaches fall short in capturing the complex temporal and structural dependencies of network traffic [21], [43].

- **State-of-the-Art AI Methods:**

AI-based methods dominate the field of synthetic data generation, particularly in recent years. Techniques such as GANs, VAEs, and Diffusion Models have emerged as the leading approaches.

- GANs (Generative Adversarial Networks): GAN-based models like CTGAN [32], WGAN with gradient penalty [16], [44], CASTGAN [35], GANBLR++ [34], [38] and have been widely used for synthetic tabular data generation. For network traffic-specific datasets, models such as FlowGAN [4], STAN [3], NetShare [6], DoppelGANger [45], and 5GT-GAN [46], MirageNet [5], have shown significant promise.
- VAEs (Variational Autoencoders): VAEs, particularly TVAE (Tabular VAE), are used as alternatives to

TABLE I  
COMPARISON OF SYNTHETIC NETWORK TRAFFIC DATA GENERATION METHODS

Category	Method	Key Features	Strengths	Weaknesses
Non-AI (Statistical)	Random Oversampling (ROS) [24]	Randomly duplicates minority class samples	Simple, low computational cost	Can lead to overfitting
	Synthetic Minority Oversampling Technique (SMOTE) [25]	Generates synthetic samples via interpolation	Generalized decision boundaries for classifiers and reduce overfitting	Doesn't consider majority class
	Adaptive Synthetic Sampling (ADASYN) [26]	Adaptive sample generation focused on hard-to-learn examples	Addresses class imbalance dynamically	Can introduce noisy samples
	Cluster Centroids (CC) [27]	Reduces majority class by replacing samples with centroids	Avoids data duplication	May oversimplify and loss key information about data structure
	Gaussian Mixture Models (GMM) [28]	Models data distribution probabilistically	Captures complex distributions	Sensitive to parameter tuning
AI-Based (Classical + Generative)	Bayesian Networks [29]	Probabilistic graphical models	Captures dependencies between features	Requires expert knowledge for structure learning
	Tabular Variational Autoencoder (TVAE) [30]	Encode-decoder architecture for learning latent representations	Learns complex distributions	Computationally expensive
	Tabular Diffusion Model (TabDDPM) [31]	Uses iterative denoising to generate data	Preserves feature dependencies	Computationally intensive
	Conditional Tabular GAN (CTGAN) [32]	Conditional generator and mode-specific normalization	Handles mixed data types effectively	Requires extensive hyperparameter tuning
	CopulaGAN [33]	Uses copula transformations in GANs	Preserves feature dependencies	Less effective in high-dimensional data
	Naive Bayes and Logistic Regression GAN (GANBLR++) [34]	Enhances GAN training with Bayesian learning	Generates both categorical and numerical data using a Dirichlet Mixture Mode	Requires accurate numerical indices
	Cascaded Tabular GAN (CasTGAN) [35]	Cascade training mechanism for tabular data	Preserves feature dependencies	Computationally intensive

GANs for tabular data generation. Recent studies, such as XIDINTFL-VAE, have demonstrated the application of VAEs in handling class imbalance within network traffic datasets [37].

- Diffusion Models: These models, including TABDDPM [31], NetDiffus [47], and NetDiffusion [7], represent an emerging paradigm in synthetic data generation. The recent NetDiffus framework has outperformed GAN-based methods, providing a 66.4% increase in the fidelity of generated data and an 18.1% increase in downstream ML tasks. They have been reported to outperform GANs and VAEs in terms of data fidelity and diversity.
- Simulators and Transformer Models: A recent study examined the feasibility of state space models, specifically Mamba, for generating packet-level synthetic network traffic data [48]. They compared the simulator with the state-of-the-art AI models, including a transformer-based TrafficGPT [49].

While these methods have individually shown great potential, they often lack direct comparative analysis under standardized evaluation metrics, making it difficult to identify the most suitable approach for specific network traffic data and its applications. Thus, a significant limitation in the current body of research is the absence of direct comparative studies evaluating Non-AI and AI (conventional and generative/SOTA) methods under unified experimental settings and standardized evaluation scenarios.

From the reviewed literature, the following key research gaps are identified:

- Limited comparative analysis of non-AI and AI (conventional and generative/SOTA) methods.
- Lack of standard evaluation metrics to assess synthetic data quality across different approaches.
- Class imbalance issues remain unexplored in the context

of network traffic datasets.

This study aims to address these gaps by systematically comparing non-AI and AI (conventional and generative/SOTA) methods using standardized evaluation metrics of fidelity and utility, with a specific focus on handling class imbalance in synthetic network traffic datasets. Therefore, for the comprehensive comparative study, we have selected twelve prominent synthetic tabular data methods. Technical details of these methods are in the next section.

### III. SYNTHETIC NETWORK TRAFFIC DATA GENERATION METHODS

The generation of synthetic network traffic data can be broadly categorized into Non-AI (statistical) methods and AI-based methods, which include both classical AI and generative AI approaches (Figure 1). This section provides a structured overview of these methods and explains the comparative evaluation approach used in this study.

#### A. Overview of Synthetic Data Generation Methods

**Non-AI Methods (Statistical Approaches):** These methods rely on mathematical models to generate synthetic data by manipulating existing samples or approximating underlying distributions. They are widely used due to their simplicity and low computational cost but often struggle to capture complex traffic dependencies [22], [50].

**AI-Based Methods (Classical AI and Generative AI):** AI-driven approaches aim to model network traffic more effectively by learning patterns from real data. Classical AI methods use probabilistic models and neural architectures to approximate traffic distributions, while generative AI techniques, such as GANs and diffusion models, focus on generating highly realistic synthetic data through adversarial training or iterative refinements [51], [52].

A summary of these methods, including their key features, strengths, and limitations, is presented in Table I. These methods were selected based on their prevalence in prior research and their relevance to the data structure of network traffic. This study focuses on comparing these approaches to determine which method best preserves fidelity, maintains class balance, and remains computationally efficient in generating synthetic network traffic data.

### B. Comparative Evaluation Approach

Rather than solely describing these methods, we systematically compare them based on the following key evaluation criteria:

- **Fidelity:** The degree to which synthetic data retains real network traffic’s statistical and temporal characteristics.
- **Utility:** The effectiveness of synthetic data in training machine learning models for network traffic classification.
- **Class Balance:** The ability of each method to generate data that reflects the true distribution of network events, addressing imbalanced datasets.
- **Scalability:** The computational efficiency of each method and its feasibility for large-scale synthetic data generation.

Each method is applied to the NSL-KDD and CIC-IDS2017 datasets to ensure an objective and standardized comparison [22], [53]. Performance is evaluated using quantitative metrics aligned with these criteria, and the results are analyzed in later sections.

This structured comparison allows us to determine the trade-offs between different approaches and provide practical recommendations for selecting the most suitable synthetic data generation techniques based on specific application needs.

## IV. EVALUATION METRICS

The following section details how each evaluation metric was measured to assess the experimental results.

- 1) **Statistical Similarity (Fidelity)** refers to the necessity for synthetic data to resemble the underlying statistical properties of real data closely [54]. This metric’s core mathematical definition is that each variable’s probability distribution in the synthetic dataset should closely match that of the corresponding variable in the real dataset [55]. However, beyond replicating the individual behavior of each variable, it is equally important to examine the interdependencies and relationships between variables. To evaluate these relationships, we compare the correlations among variables in both datasets [56]. Another crucial aspect is the preservation of the data structure. For instance, if a variable is binary in the real data, it must remain binary in the synthetic data to ensure logical consistency. Similarly, suppose a variable is defined to have values greater than or equal to zero (such as packet size in network traffic data). In that case, the corresponding variable in the synthetic data should also respect this constraint. Failing to maintain these structural properties can lead to synthetic data that,

while statistically similar on a superficial level, fails to capture the logical and contextual nuances necessary for accurate downstream analyses and applications [57]. In this study, the following metrics were used to measure the statistical similarity (fidelity), as shown in results Table VI and Table VII:

### A. Experiment Setup

- a) **Data Structure (DS):** This metric checks whether the synthetic data adheres to the original data’s logical minimum and maximum values. For example, a binary column should contain only binary values, ensuring no out-of-bound entries.

Table VI and Table VII include a column where each value is calculated based on whether all binary and boolean variables adhere to their expected values. This means that binary variables should be either 0 or 1, and boolean variables should be either `TRUE` or `FALSE`. If the synthetic data generated by a method does not meet these structural requirements, the result will be marked as `NO`.

For example, in Table VI, the first synthetic data set was generated by ROS, and the value in the DS column is `YES`, indicating that all binary and boolean variables adhered to their expected values. However, the next row, corresponding to SMOTE, shows `NO` under the DS column because, apart from the target variable in the NSL-KDD data, the binary and boolean variables did not conform to their expected values. This evaluation method helps assess whether each approach maintains the required data boundaries.

- b) **Correlation (Corr):** This assessment is performed by examining the correlation heatmap, which represents the absolute values of the correlation coefficients for all variables (see Figure 2 and Figure 3 for key NSL KDD heatmaps). It evaluates the strength and direction of relationships in both real and synthetic datasets. If there is a significant difference between the absolute correlation heatmap of the synthetic data and the real data, as well as in the heatmap of the absolute correlation differences, the result is marked as `NO`. Otherwise, it is marked as `YES` in the Corr column of all of the results tables.
- c) **Probability Distribution (PD):** This compares the underlying probability distribution of each variable between the real and synthetic datasets, ensuring that their statistical properties are aligned. To calculate the statistics for the PD comparison, the following steps were undertaken:
  - i) For every variable in the synthetic and real datasets, the underlying PD was estimated using Kernel Density Estimation (KDE).

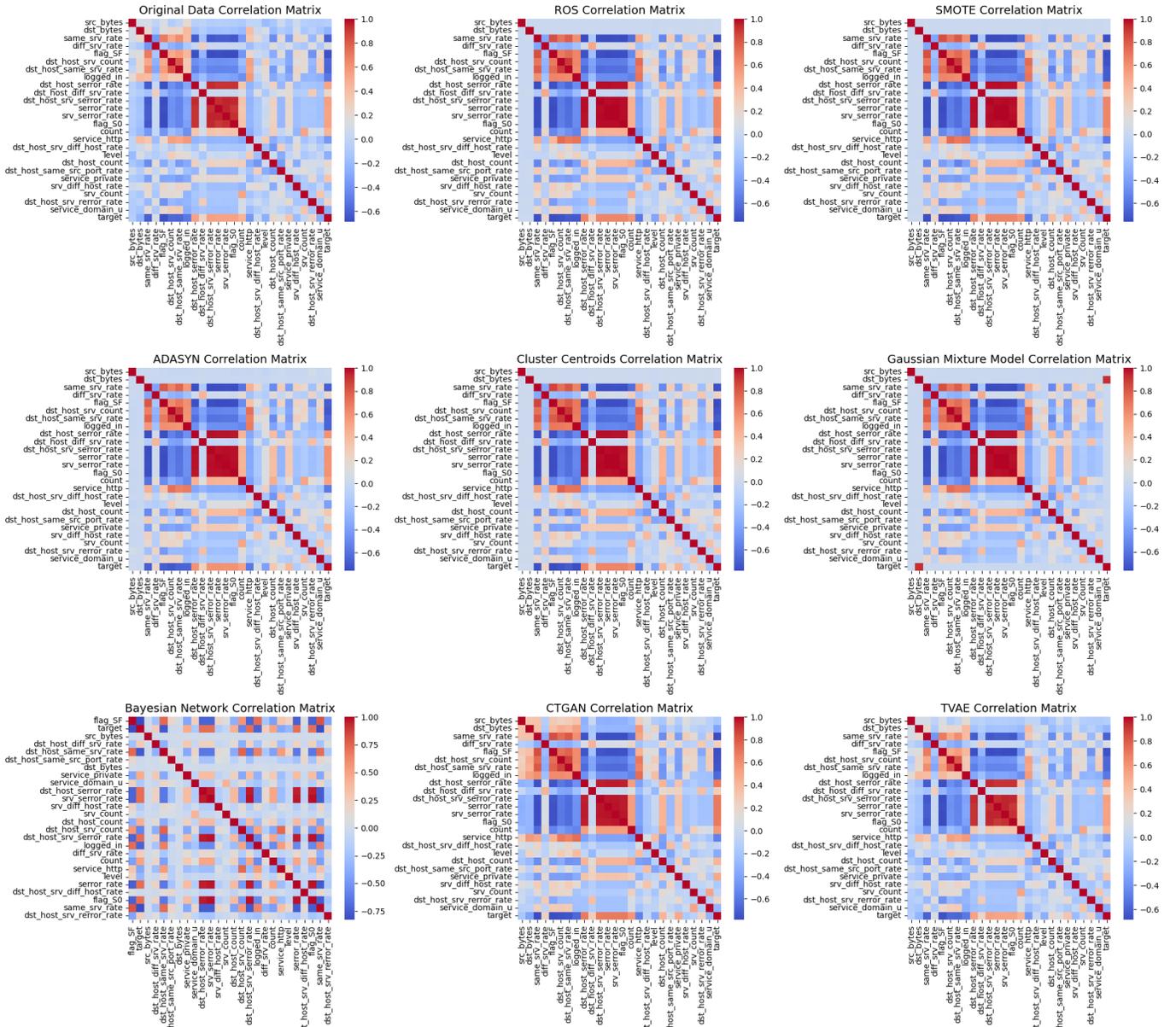


Fig. 2. Correlation Heatmap for NSL KDD (Part 1/2)

- ii) Each variable's PD was visualized. For clarity, individual figures showing the real and synthetic PDs were generated. In some cases, side-by-side figures were created to provide a direct comparison (Figure 2). Complete PD comparison graphs for both NSL-KDD and CIC-IDS2017 are provided in Appendix A.
- iii) The alignment of the PD was quantified using the following formula:

$$PD(\%) = \frac{\text{Number of Variables with Different PD}}{\text{Total Number of Variables}} \times 100 \quad (1)$$

where "Number of Variables with Different

PD indicates the count of variables whose PD differed between the real and synthetic datasets.

For the NSL-KDD dataset, which contains 26 total variables, the results are summarized in Table VI. Under the PD column, for the ROS method, all variables exhibited identical PDs according to distribution curves in Figure 4, resulting in a PD of 0% difference. Conversely, for the Adaptive Synthetic Sampling method, one variable out of 26 exhibited a different PD (see Appendix A for a complete graphical PD comparison). Using Equa-

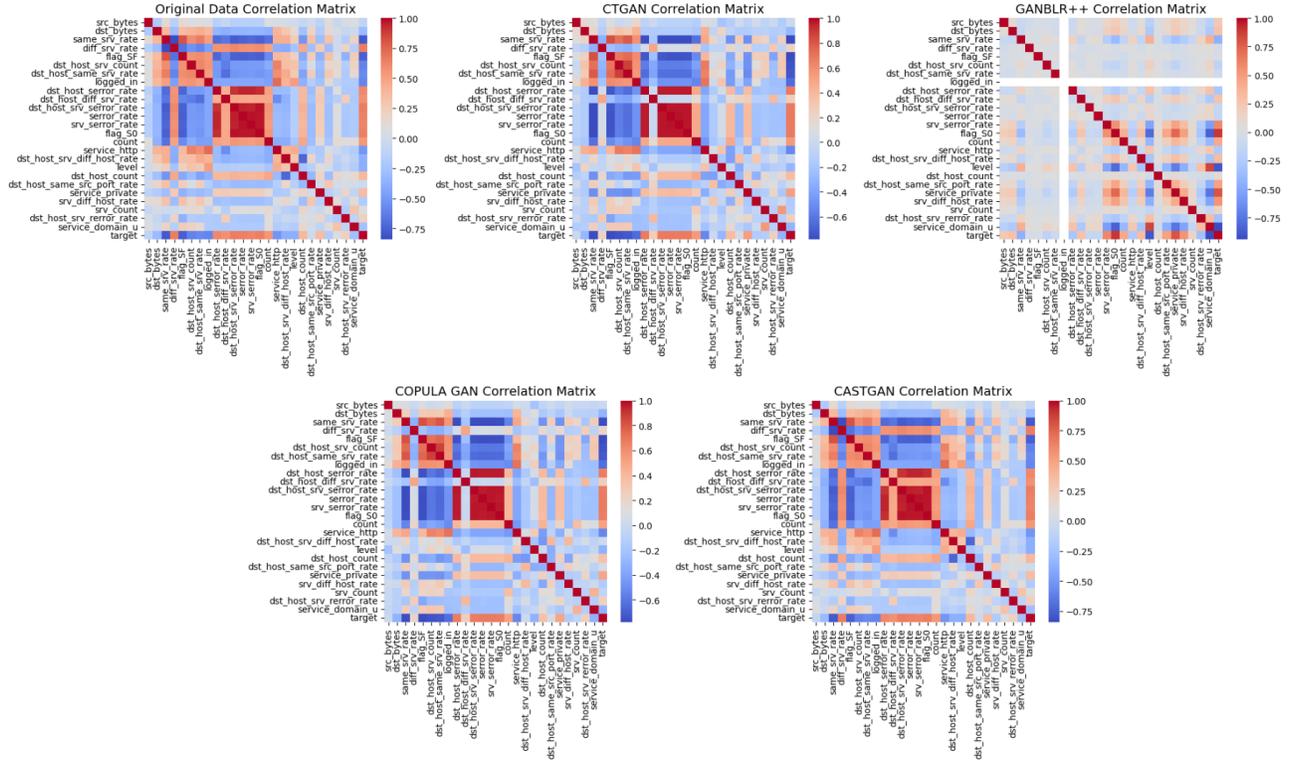


Fig. 3. Correlation Heatmap for NSL KDD (Part 2/2)

tion 1, the percentage difference was calculated as:

$$PD (\%) = \frac{1}{26} \times 100 \approx 3.8\% \quad (2)$$

- 2) **Performance (utility/accuracy)** is a widely used metric to evaluate the utility of synthetic data, particularly in the context of ML and DL models [59]. This metric compares how models trained on synthetic data perform relative to those trained on real data [60]. The most common approach is to contrast TRTR (train and test on real data) with TSTR (train on synthetic data and test on real data). By assessing the model’s accuracy, F1 score, precision, and recall, researchers can determine the efficacy of the synthetic data. Suppose the model’s performance metrics in the TSTR scenario are comparable to or exceed those in the TRTR scenario. In that case, the synthetic data is considered a viable alternative to the real data [38].

Moreover, a well-performing TSTR indicates that the synthetic data captures the essential patterns and relationships within the real data, making it suitable for model development, validation, and deployment in privacy-sensitive environments. However, it’s crucial to note that synthetic data’s utility is not solely dependent on performance metrics; it also relies on the data’s ability to generalize across different models and tasks, which should be considered in comprehensive evaluations [61].

Thus, in this study, the accuracy of TRTR was compared with TSTR for measuring ML utility, as reported in Table VI and Table VII.

- 3) **Class Balance (CB)** is a critical metric for evaluating the quality of synthetic tabular data, particularly in classification tasks. Class balance refers to the distribution of instances across different classes in the dataset [62]. Maintaining a balanced class distribution in real-world datasets is essential for training robust ML models, as imbalanced courses can lead to biased models that perform poorly on underrepresented classes [63].

$$\text{Class Balance Difference } (\%) = \frac{|P_{\text{synthetic,Normal}} - P_{\text{synthetic,Attack}}|}{\times 100} \quad (3)$$

Thus, the CB results mentioned in Table VI and Table VII represent how much difference exists between the NORMAL and ATTACK cases. If the classes in the synthetic data are the same, there is zero difference; otherwise, the difference percentage is calculated based on the equation. This equation 3 gives 0% diff if the normal% and attack% classes are equal, i.e., follow a 50-50 ratio in the whole data. Thus, in Table VI, under the CB column for ROS and SMOTE, it is mentioned 0% diff, and for the Adaptive Synthetic Sampling, it is 0.14% diff.

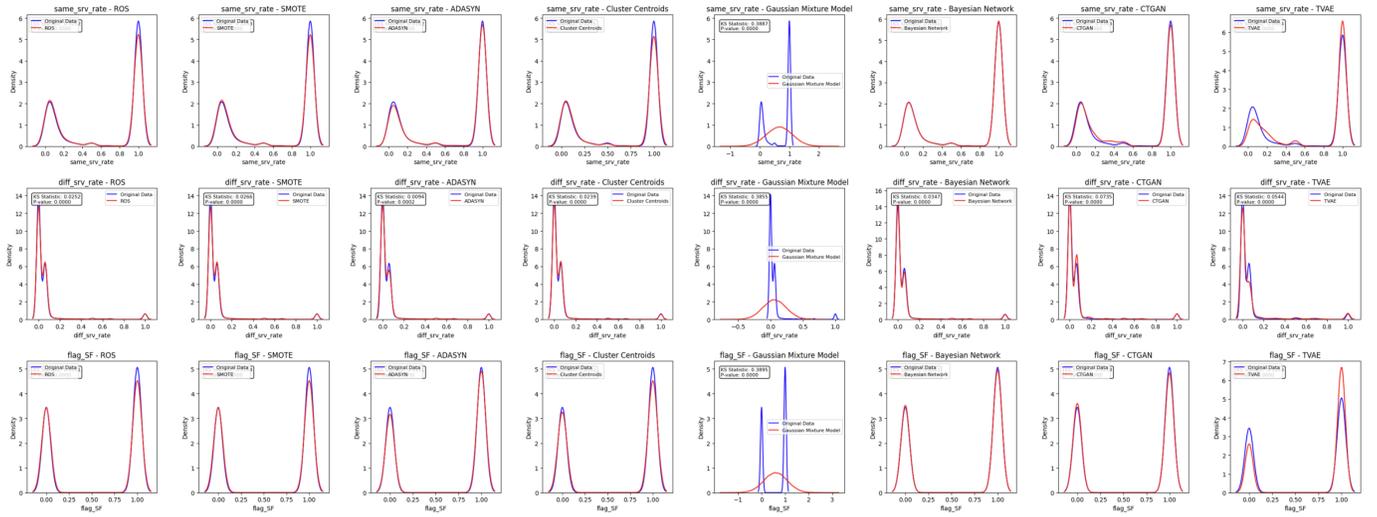


Fig. 4. Probability Distribution Comparison for NSL KDD.

## V. EXPERIMENTAL SETUP

### A. Mutual Information

Mutual Information (MI) measures the mutual dependence between two random variables (linear or non-linear), initially introduced by Shannon in his foundational work on information theory [64]. It quantifies the amount of information gained about one random variable through the observation of another [65]. Unlike the correlation coefficient, limited to real-valued random variables and linear relationships, MI is more general, capturing both linear and non-linear dependencies [66]. The concept of MI is closely related to entropy, a fundamental notion in information theory that measures the expected amount of information in a random variable [67]. MI can be understood as the expected value of the pointwise mutual information (PMI), reflecting how much the joint distribution of two variables deviates from the product of their marginal distributions [68].

1) *Comparison with Other Methods:* The rationale for selecting MI over other methods lies in its comprehensive mathematical framework. Initially, this study considered two general approaches for choosing the best features based on their dependence on the target feature and their relationships. Correlation was the first and most obvious choice, mathematically defined as the Pearson correlation coefficient (equ: 4). However, correlation only accounts for linear relationships [69]. In real-world scenarios, non-linear relationships often need to be identified and quantified to select the optimal set of features.

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \quad (4)$$

where  $\text{Cov}(X, Y)$  is the covariance between variables  $X$  and  $Y$ , and  $\sigma_X$  and  $\sigma_Y$  are the standard deviations of  $X$  and  $Y$ , respectively. While correlation is effective for capturing linear relationships, it does not account for the non-linear

relationships that often exist in real-world data, limiting its effectiveness for feature selection.

The second option was a tree-based AI technique, known for its robustness and ability to identify essential features rigorously [70]. However, tree-based methods suffer from the "black box" problem, offering little interpretability regarding why certain features are selected.

While approaches like Information Gain (IG) have been used for feature selection, as demonstrated by one study [71], Information Gain only measures the relationship between individual features and the target variable, overlooking interdependencies between features [72]. Mathematically, it measures the reduction in entropy when a feature  $X$  is used to predict the target variable  $Y$ .

$$IG(Y, X) = H(Y) - H(Y|X) \quad (5)$$

where  $H(Y)$  is the entropy of the target variable  $Y$ , and  $H(Y|X)$  is the conditional entropy of  $Y$  given  $X$ . Although Information Gain quantifies the relationship between individual features and the target variable, it ignores interdependencies among the features.

In contrast, MI captures the relationship between features and the target and the dependencies among the features themselves, providing a more holistic evaluation [73]. MI is defined as:

$$MI(X, Y) = H(X) + H(Y) - H(X, Y) \quad (6)$$

where  $H(X)$  and  $H(Y)$  are the entropies of variables  $X$  and  $Y$ , and  $H(X, Y)$  is their joint entropy. By capturing non-linear dependencies and interactions among features, MI offers a more comprehensive framework for feature selection.

Compared to other AI-based feature selection methods, such as recursive feature elimination or embedded methods based on decision trees [74], MI has the upper hand by

TABLE II  
SUMMARY OF 26 ENCODED SELECTED FEATURES FROM NSL-KDD DATASET [58]

S.No	Feature	Description
1	src_bytes	Number of data bytes from source to destination
2	dst_bytes	Number of data bytes from destination to source
3	same_srv_rate	Percentage of connections to the same service
4	diff_srv_rate	Percentage of connections to different services
5	flag_SF	Connection status with a normal connection (SF: "Normal")
6	dst_host_srv_count	Number of connections to the same service as the current connection in the past 100 connections
7	dst_host_same_srv_rate	Percentage of connections to the same service for a destination host
8	logged_in	1 if successfully logged in; 0 otherwise
9	dst_host_serror_rate	Percentage of connections that have "SYN" errors
10	dst_host_diff_srv_rate	Percentage of connections to different services for a destination host
11	dst_host_srv_serror_rate	Percentage of connections that have "SYN" errors for a destination host
12	serror_rate	Percentage of connections that have "SYN" errors
13	srv_serror_rate	Percentage of connections that have "SYN" errors for the same service
14	flag_S0	Connection status where no data packets were exchanged (S0: "No Data Exchange")
15	count	Number of connections to the same host as the current connection in the past 2 seconds
16	service_http	HTTP service (1 if used, 0 otherwise)
17	dst_host_srv_diff_host_rate	Percentage of connections to different hosts on the same service
18	level	Threat level of the connection
19	dst_host_count	Number of connections to the same destination host in the past 100 connections
20	dst_host_same_src_port_rate	Percentage of connections with the same source port to the destination host
21	service_private	Private network service (1 if used, 0 otherwise)
22	srv_diff_host_rate	Percentage of connections to different hosts for the same service
23	srv_count	Number of connections to the same service as the current connection in the past 2 seconds
24	dst_host_srv_error_rate	Percentage of connections that have "REJ" errors for a destination host
25	service_domain_u	Domain name service (DNS) (1 if used, 0 otherwise)
26	target	Class label indicating if the connection is normal or an attack

being model-agnostic. Many AI-based techniques are tied to specific algorithms or models, which may introduce biases or limit generalizability across different datasets or classifiers. MI, however, evaluates feature importance based on statistical dependencies independent of any specific model, making it a versatile and unbiased approach [75]. This allows MI to provide a more robust selection process, especially in scenarios where feature relationships are complex and nonlinear, a common occurrence in network data.

2) *Application of Mutual Information in this Study*: In this study, MI is crucial for measuring the dataset's dependency between different features (variables). It helps understand the strength of relationships between variables, which is essential for modeling and analysis. For example, mutual information was used to quantify the relation between each feature  $x_i$  and the target  $y$  and between each feature to identify interdependencies within the features. This step was crucial to filter down the essential features that explain the target variables and to cluster the features exhibiting similar behavior.

The mutual information score of every feature was determined concerning the target, and the top 25% of the features were selected based on the mutual information weights (Table IV & Table V). This way, the computation was less intensive, the information was high, and the results were more inter-

pretable.

## B. Datasets & Pre-processing

- 1) **NSL-KDD** dataset, a refined version of the original KDD Cup 1999 dataset, is widely used in cybersecurity to evaluate IDS [58]. It addresses some critical issues present in its predecessor, such as redundant records and imbalanced distribution, thereby providing a more accurate and reliable benchmark for performance evaluation. For this study, the training portion of the NSL-KDD dataset, originally consisting of 125,973 instances and 42 columns (41 features and 1 (binary) target), was used. The preprocessing of the NSL-KDD dataset involved several essential steps to prepare it for synthetic data generation and evaluation. Initially, categorical columns were encoded to transform non-numeric data into a suitable format for ML algorithms. This encoding process ensured that all categorical variables were converted into numerical values, facilitating the subsequent analysis. Following the encoding, feature extraction was performed to identify the most significant features using Mutual Information. This process selected the features' top 25% (Q1) based on the mutual information weights. This step resulted in a reduced feature set of 26 encoded

TABLE III  
SUMMARY OF 21 SELECTED FEATURES FROM CIC-IDS2017 DATASET [12]

S.No	Feature	Description
1	Average Packet Size	Average size of the packets in the flow
2	Packet Length Std	Standard deviation of the packet lengths in the flow
3	Packet Length Variance	Variance of the packet lengths in the flow
4	Packet Length Mean	Mean length of the packets in the flow
5	Total Length of Bwd Packets	Total number of bytes in backward (Bwd) packets
6	Subflow Bwd Bytes	Number of bytes in the backward sub-flow
7	Destination Port	Port number of the destination host
8	Avg Bwd Segment Size	Average size of backward segments
9	Bwd Packet Length Mean	Mean length of backward packets
10	Init_Win_bytes_forward	Initial window size in bytes for forward direction
11	Subflow Fwd Bytes	Number of bytes in the forward sub-flow
12	Total Length of Fwd Packets	Total number of bytes in forward (Fwd) packets
13	Max Packet Length	Maximum length of a packet in the flow
14	Bwd Packet Length Max	Maximum length of a backward packet
15	Init_Win_bytes_backward	Initial window size in bytes for backward direction
16	Fwd Packet Length Max	Maximum length of a forward packet
17	Fwd Packet Length Mean	Mean length of forward packets
18	Avg Fwd Segment Size	Average size of forward segments
19	Flow IAT Max	Maximum time interval between packets in the flow
20	Flow Bytes/s	Rate of flow in bytes per second
21	target	Class label indicating if the flow is benign or malicious

features (Table II), which were the most important (based on the information weights; Table IV) for accurately representing the data. This refined feature set was then used in the data generation phase, ensuring that the synthetic data closely mirrored the critical characteristics of the original dataset.

- 2) **CIC-IDS2017** dataset, part of the Canadian Institute for Cybersecurity’s intrusion detection dataset collection, is widely employed for evaluating cybersecurity systems [12]. It provides a comprehensive network traffic data set with diverse attack types and normal traffic, offering a robust benchmark for performance assessment. CIC-IDS2017 has one week of captured network traffic data in eight comma-separated files (CSV). These files are generated based on the attack types. There are 2,830,743 instances of all eight files and 79 columns (78 numerical features and one binary target).

Preprocessing of the CIC-IDS2017 dataset involved several vital steps to ready it for synthetic data generation and evaluation. First, all the files were individually tested for null values and duplicated columns. In this process, it was found that there were some infinity values in the data. Thus, all the null values, infinity, and duplicated column “*Fwd Header Length*” were removed. Secondly, all the cleaned data files were concatenated into a single CSV file for further feature selection. Lastly, mutual information was performed to select the variables’ top 25% (Q1) based on the mutual information weights. This step resulted in a reduced feature set of 21 features (Table III), which were determined to be the most important (based on the information weights; V) for accurately representing the data. This reduced feature set was then employed in the synthetic data generation phase, ensuring that the synthetic data accurately reflected the

essential characteristics of the original dataset.

TABLE IV  
MUTUAL INFORMATION SCORE OF NSL-KDD FEATURES

Feature	Mutual Information Score
src_bytes	0.566864
dst_bytes	0.439281
same_srv_rate	0.369288
diff_srv_rate	0.361895
flag_SF	0.341828
dst_host_srv_count	0.335993
dst_host_same_srv_rate	0.309832
logged_in	0.292075
dst_host_serror_rate	0.286589
dst_host_diff_srv_rate	0.283874
dst_host_srv_serror_rate	0.281332
serror_rate	0.278666
srv_serror_rate	0.269186
flag_S0	0.263399
count	0.262733
service_http	0.191343
dst_host_srv_diff_host_rate	0.189822
level	0.153819
dst_host_count	0.144479
dst_host_same_src_port_rate	0.131316
service_private	0.118493
srv_diff_host_rate	0.099441
srv_count	0.062476
dst_host_srv_rerror_rate	0.062244
service_domain_u	0.048430

1) *Hardware Setup*: The experiments were performed on a workstation with a 13th Gen Intel® Core™ i9-13900 processor. This processor operates at 2000 MHz and features 24 cores with 32 logical processors. The system is equipped with 32 GB of RAM. A dedicated GPU, NVIDIA GeForce RTX 4090, was used for these experiments.

2) *Software Setup*: The software environment for the experiments was based on Microsoft Windows 11 Pro. The primary programming language used was Python 3.12.4, which

TABLE V  
MUTUAL INFORMATION SCORE OF CIC-IDS2017 FEATURES

Feature	Mutual Information Score
Average Packet Size	0.347112
Packet Length Std	0.342188
Packet Length Variance	0.342019
Packet Length Mean	0.319635
Total Length of Bwd Packets	0.296955
Subflow Bwd Bytes	0.296882
Destination Port	0.291245
Avg Bwd Segment Size	0.287676
Bwd Packet Length Mean	0.287483
Init_Win_bytes_forward	0.287174
Subflow Fwd Bytes	0.284528
Total Length of Fwd Packets	0.284264
Max Packet Length	0.264060
Bwd Packet Length Max	0.263210
Init_Win_bytes_backward	0.250188
Fwd Packet Length Max	0.246537
Fwd Packet Length Mean	0.213223
Avg Fwd Segment Size	0.213065
Flow IAT Max	0.212817
Flow Bytes/s	0.209784

provides a stable and versatile platform for implementing various data science and ML techniques. Essential libraries and frameworks include NumPy for numerical operations, Pandas for data manipulation, Seaborn and matplotlib for data visualization, and sci-kit-learn for ML utilities. The experiments also leveraged PyTorch for deep learning tasks. For handling class imbalance, the imbalanced-learn library (learn) was employed. The libraries are crucial for the GAN models, including SDV (Synthetic Data Vault) for tabular data generation, as they support models like CTGAN, TVAE, and CopulaGAN. Due to the lack of pre-existing packages, their official GitHub repositories were also utilized for specialized GAN models such as GANBLR and CASTGAN. PGMPy and SciPy were also used for probabilistic graphical models and scientific computations.

## VI. RESULTS AND DISCUSSION

This section presents an in-depth comparative evaluation of twelve synthetic data generation methods applied to the NSL-KDD and CIC-IDS2017 datasets. The analysis is structured around four key objectives: fidelity, utility, class balance, and computational efficiency. Tables VI and VII summarize these results for both datasets.

The evaluation focuses on two distinct categories of methods: Non-AI (Statistical) and AI-based (Classical + Generative AI). The performance of these methods is assessed based on statistical similarity (SS), class balance (CB), and machine learning utility (measured as accuracy in TRTR vs. TSTR). Additionally, AI-based methods' computational cost and training efficiency are critically analyzed.

### A. Statistical Methods

Statistical methods such as ROS, SMOTE, ADASYN, and Cluster Centroids (CC) are widely used for addressing class imbalance in network traffic datasets. However, these methods primarily rely on resampling strategies rather than truly

modeling the underlying data distribution. While SMOTE and ADASYN generate synthetic samples, they do so through interpolation, which may fail to capture complex feature dependencies. Additionally, statistical oversampling techniques can lead to overfitting in downstream machine learning tasks, as the generated samples lack inherent variability compared to real network traffic. Prior studies [X, Y] indicate that such methods are effective in reducing bias in classification but perform poorly when high-dimensional feature interactions are critical for network anomaly detection. They exhibit excellent class balance maintenance, achieving 0% CB difference in most cases. However, these methods primarily resample existing data points rather than generating truly novel synthetic data, which limits their effectiveness in capturing complex relationships between features.

- SMOTE and CC outperform ROS by enhancing ML utility while preserving statistical similarity.
- GMM performs poorly due to its parametric limitations, struggling to model high-dimensional network traffic distributions (99.1% PD diff in NSL-KDD, 61.9% in CIC-IDS2017).
- Despite their fidelity limitations, statistical methods achieve near-perfect accuracy (0.999) in TRTR and TSTR settings, showing they can still be useful for ML training in balanced data scenarios.

### B. AI-Based Methods

AI-based methods exhibit stronger generative capabilities but face trade-offs in fidelity, utility, and computational cost:

- CTGAN and CopulaGAN achieve the best balance, with high fidelity (Corr: Yes, PD diff: 7.7% (NSL-KDD), Corr: Yes, PD diff: 0% (CIC-IDS2017)) and strong utility (accuracy: 0.9667 (NSL-KDD) and 0.9538 (CIC-IDS2017)).
- TVAE is overall highly effective in NSL-KDD but struggles with class balance in CIC-IDS2017 (CB diff: 69.4%), indicating possible limitations in handling imbalanced network traffic distributions.
- Diffusion models (TABDDPM) perform poorly, particularly in high-dimensional tabular data, due to their reliance on iterative denoising rather than structured feature learning (PD diff: 98.3% (NSL-KDD)).

### C. Computational Cost and Training Efficiency

The computational efficiency of AI-based methods varies significantly, impacting their scalability for large-scale network traffic simulation:

- Bayesian Networks (BN) failed on CIC-IDS2017 (Table VII) due to excessive memory requirements when constructing initial feature relationships. This highlights the challenge of applying probabilistic graphical models to large network traffic datasets.
- TABDDPM could not generate results for CIC-IDS2017 (Table VII), indicating prohibitively high computational costs. Diffusion-based models, while powerful for image data, struggle with tabular datasets due to their need for extensive forward and reverse sampling processes.

TABLE VI  
COMPARISON OF SYNTHETIC DATA ON NSL-KDD DATASET

Category	Method	SS			CB	Accuracy
		DS	Corr	PD		
Non-AI (Statistical)	ROS	Yes	Yes	0% diff	0% diff	0.9998
	SMOTE	No	Yes	0% diff	0% diff	0.9999
	ADASYN	No	Yes	3.8% diff	0.14% diff	0.9998
	CC	No	Yes	0% diff	0% diff	0.9995
	GMM	No	Yes	99.1% diff	99.9% diff	0.5314
AI-Based (Classical + Generative)	BN	No	No	0% diff	17.2% diff	0.9587
	TVAE	Yes	Yes	23.1% diff	26.7% diff	0.9807
	TABDDPM	No	No	98.3% diff	34.2% diff	0.51
	CTGAN	Yes	Yes	7.7% diff	6.7% diff	0.9667
	GANBLR++	No	No	99% diff	8.8% diff	0.5916
	CopulaGAN	Yes	Yes	7.7% diff	5.4% diff	0.9761
	CasTGAN	No	Yes	27% diff	11.1% diff	0.9582

"SS" = Statistical Similarity, "DS" = Data Structure, "Corr" = Correlation, "PD" = Probability Distribution, "CB" = Class Balance, "Accuracy" = Machine Learning Utility on TRTR vs TSTR, "TRTR" = Train and Test on Real Data = **0.9995**, "TSTR" = Train on Synthetic and Test on Real Data

TABLE VII  
COMPARISON OF SYNTHETIC DATA ON CIC-IDS2017 DATASET

Category	Method	SS			CB	Accuracy
		DS	Corr	PD		
Non-AI (Statistical)	ROS	Yes	Yes	0% diff	0% diff	0.9991
	SMOTE	Yes	Yes	0% diff	0% diff	0.9990
	ADASYN	Yes	Yes	0% diff	0.2% diff	0.9986
	CC	No	Yes	0% diff	0% diff	0.9930
	GMM	No	No	61.9% diff	46.1% diff	0.4509
AI-Based (Classical + Generative)	BN	-	-	-	-	-
	TVAE	Yes	No	1.4% diff	69.4% diff	0.9718
	TABDDPM	-	-	-	-	-
	CTGAN	Yes	Yes	0% diff	5% diff	0.95381
	GANBLR++	No	No	47.6% diff	64.4% diff	0.5378
	CopulaGAN	Yes	Yes	0% diff	5% diff	0.9755
	CasTGAN	No	No	4.8% diff	58.3% diff	0.9734

"SS" = Statistical Similarity, "DS" = Data Structure, "Corr" = Correlation, "PD" = Probability Distribution, "CB" = Class Balance, "Accuracy" = Machine Learning Utility on TRTR vs TSTR, "TRTR" = Train and Test on Real Data = **0.9995**, "TSTR" = Train on Synthetic and Test on Real Data

- CTGAN and CopulaGAN balance computational efficiency and fidelity, making them the most practical for synthetic network traffic data where scalability is critical.

Key technical inference from the abovementioned results is that GAN-based models, such as CTGAN and CopulaGAN, perform exceptionally well because they utilize mode-specific normalization and conditional generation, which helps preserve the feature dependencies crucial for generating realistic network traffic samples. This ability to maintain the relationships between features contributes to their effectiveness in generating high-quality synthetic data [32], [76].

On the other hand, diffusion models face significant challenges, particularly when applied to large tabular datasets. These models rely on iterative denoising, a process that becomes computationally expensive as the dataset size increases. Unlike image data, which has spatial dependencies that aid the diffusion process, tabular data lacks this inherent structure, making the data generation process less efficient. Furthermore, diffusion models struggle with discrete categorical features, which are common in network traffic datasets, further hindering their performance [77].

Gaussian Mixture Models also fail to capture the complex

nature of network traffic data. These models assume that the data follows underlying Gaussian distributions, but they are unable to account for the multimodal and varied traffic behaviors typically found in network traffic datasets [78].

For practical recommendations, techniques like SMOTE and ROS are still viable options for addressing class imbalance in network traffic datasets. When it comes to high-fidelity data generation, CTGAN and CopulaGAN are the top performers, surpassing other models in generating realistic and useful synthetic data. However, Bayesian Networks and Diffusion Models should be avoided for large-scale datasets due to their high computational costs and inefficiencies in handling such data.

## VII. FUTURE WORK

Building upon the insights derived from this study, potential directions for future investigations include:

- 1) **Real vs AI class balance:** A predominant challenge in ML tasks is ensuring balanced classes for unbiased predictions. This requirement raises questions about the real-world balance of attack classes, necessitating thorough examination. It is crucial to scrutinize whether

achieving balance among ML classes facilitates accurate data synthesis and optimizes ML tasks.

- 2) **Benchmarking the boundary of attack and normal traffic:** This research applies generative AI techniques to analyze the distinction between normal and abnormal network traffic rigorously. Through this exploration, we aim to enhance our comprehension of threshold dynamics, significantly contributing to developing more resilient IDS.
- 3) **Fidelity vs Privacy:** Enhancing the fidelity and utility of synthetic data often poses privacy risks. Thus, it is imperative to identify an optimized balance between fidelity and privacy, particularly concerning cybersecurity network traffic datasets. This endeavor requires careful consideration to safeguard individual privacy while maintaining the data's utility.
- 4) **Data dynamic features:** Develop and refine feature selection methods targeting data dynamic features in network flows. These techniques should adapt to evolving network environments, identifying the most relevant features for real-time IDS.
- 5) **Adaptive IDS:** Design IDS that dynamically adjusts to changes in network traffic patterns. This would involve real-time analysis and feature selection to detect new threats emerging as the network evolves.
- 6) **Loss functions:** Conduct a comparative analysis of various generative AI loss functions tailored to the study's specific context. This will help identify which loss functions are most effective in optimizing the application's model performance, providing insights into their strengths and limitations in different scenarios.
- 7) **Correlation vs. Dependency:** To examine the deeper relationships between variables by exploring logical and conditional dependencies, going beyond the surface-level insights provided by simple correlation measures. This will allow a more critical understanding of how variables interact, particularly in complex datasets where correlations may not fully capture the underlying dependencies.

## VIII. CONCLUSION

This study systematically evaluates synthetic network traffic data generation methods, comparing statistical and AI-based approaches. Our findings reveal that while statistical methods effectively handle class imbalance, they lack the ability to preserve complex feature dependencies. GAN-based models, particularly CTGAN and CopulaGAN, offer the best trade-off between fidelity and utility, making them preferable for generating realistic synthetic network traffic data. Diffusion models, despite their recent advances, exhibit significant computational challenges when applied to large-scale tabular datasets. The study highlights the need for hybrid approaches integrating statistical and AI techniques to improve class balance while maintaining data realism. Future research should explore optimizing generative AI loss functions, benchmarking real versus

AI-generated class balance, and addressing privacy concerns while enhancing synthetic data fidelity.

## IX. ACKNOWLEDGEMENTS

The work presented here was partly financed by the European Celtic+ and Swedish Vinnova project "CISSAN – Collective Intelligence Supported by Security Aware Nodes".

## REFERENCES

- [1] A. Shahraki, M. Abbasi, A. Taherkordi, and A. D. Jurcut, "A comparative study on online machine learning techniques for network traffic streams analysis," *Computer Networks*, vol. 207, p. 108836, 2022.
- [2] J. L. Guerra, C. Catania, and E. Veas, "Datasets are not enough: Challenges in labeling network traffic," *Computers & Security*, vol. 120, p. 102810, 2022.
- [3] S. Xu, M. Marwah, M. Arlitt, and N. Ramakrishnan, "Stan: Synthetic network traffic generation with generative neural models," in *Deployable Machine Learning for Security Defense: Second International Workshop, MLHat 2021, Virtual Event, August 15, 2021, Proceedings 2*, pp. 3–29, Springer, 2021.
- [4] L. D. Manocchio, S. Layeghy, and M. Portmann, "Flowgan-synthetic network flow generation using generative adversarial networks," in *2021 IEEE 24th International Conference on Computational Science and Engineering (CSE)*, pp. 168–176, IEEE, 2021.
- [5] S. K. Nukavarapu, M. Ayyat, and T. Nadeem, "Miragenet-towards a gan-based framework for synthetic network traffic generation," in *GLOBECOM 2022-2022 IEEE Global Communications Conference*, pp. 3089–3095, IEEE, 2022.
- [6] Y. Yin, Z. Lin, M. Jin, G. Fanti, and V. Sekar, "Practical gan-based synthetic ip header trace generation using netshare," in *Proceedings of the ACM SIGCOMM 2022 Conference*, pp. 458–472, 2022.
- [7] X. Jiang, S. Liu, A. Gember-Jacobson, A. N. Bhagoji, P. Schmitt, F. Bronzino, and N. Feamster, "Netdiffusion: Network data augmentation through protocol-constrained traffic generation," *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, vol. 8, no. 1, pp. 1–32, 2024.
- [8] A. Mozo, Á. González-Prieto, A. Pastor, S. Gómez-Canaval, and E. Tálavera, "Synthetic flow-based cryptomining attack generation through generative adversarial networks," *Scientific reports*, vol. 12, no. 1, p. 2091, 2022.
- [9] A. Kotal, A. Piplai, S. S. L. Chukkappalli, and A. Joshi, "Privetab: Secure and privacy-preserving sharing of tabular data," in *Proceedings of the 2022 ACM on International Workshop on Security and Privacy Analytics*, pp. 35–45, 2022.
- [10] D. Ganji and C. Chakrabortii, "Towards data generation to alleviate privacy concerns for cybersecurity applications," in *2023 IEEE 47th Annual Computers, Software, and Applications Conference (COMPSAC)*, pp. 1447–1452, IEEE, 2023.
- [11] S. Layeghy, M. Gallagher, and M. Portmann, "Benchmarking the benchmark—comparing synthetic and real-world network ids datasets," *Journal of Information Security and Applications*, vol. 80, p. 103689, 2024.
- [12] I. Sharafaldin, A. H. Lashkari, A. A. Ghorbani, *et al.*, "Toward generating a new intrusion detection dataset and intrusion traffic characterization," *ICISSp*, vol. 1, pp. 108–116, 2018.
- [13] V. Dixit, R. Selvarajan, T. Aldwairi, Y. Koshka, M. A. Novotny, T. S. Humble, M. A. Alam, and S. Kais, "Training a quantum annealing based restricted boltzmann machine on cybersecurity data," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 6, no. 3, pp. 417–428, 2021.
- [14] M. Sabet, P. Palanisamy, and S. Mishra, "Scalable modular synthetic data generation for advancing aerial autonomy," *Robotics and Autonomous Systems*, vol. 166, p. 104464, 2023.
- [15] I. H. Sarker, A. Kayes, S. Badsha, H. Alqahtani, P. Watters, and A. Ng, "Cybersecurity data science: an overview from machine learning perspective," *Journal of Big data*, vol. 7, pp. 1–29, 2020.
- [16] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative adversarial networks," in *International conference on machine learning*, pp. 214–223, PMLR, 2017.
- [17] Y. N. Rao and K. Suresh Babu, "An imbalanced generative adversarial network-based approach for network intrusion detection in an imbalanced dataset," *Sensors*, vol. 23, no. 1, p. 550, 2023.

- [18] Z. Lin, Y. Shi, and Z. Xue, "Idsgan: Generative adversarial networks for attack generation against intrusion detection," in *Pacific-asia conference on knowledge discovery and data mining*, pp. 79–91, Springer, 2022.
- [19] N. Sivaroopan, C. Madarasingha, S. Muramudalige, G. Jourjon, A. Jayasumana, and K. Thilakarathna, "Synig: Synthetic network traffic generation through time series imaging," in *2023 IEEE 48th Conference on Local Computer Networks (LCN)*, pp. 1–9, IEEE, 2023.
- [20] S. A. Alex, "Imbalanced data learning using smote and deep learning architecture with optimized features," *Neural Computing and Applications*, vol. 37, no. 2, pp. 967–984, 2025.
- [21] E. Papadaki, A. G. Vrahatis, and S. Kotsiantis, "Exploring innovative approaches to synthetic tabular data generation," *Electronics*, vol. 13, no. 10, p. 1965, 2024.
- [22] A. Figueira and B. Vaz, "Survey on synthetic data generation, evaluation methods and gans," *Mathematics*, vol. 10, no. 15, p. 2733, 2022.
- [23] A. Kiran and S. S. Kumar, "A comparative analysis of gan and vae based synthetic data generators for high dimensional, imbalanced tabular data," in *2023 2nd International Conference for Innovation in Technology (INOCON)*, pp. 1–6, IEEE, 2023.
- [24] A. S. Dina, A. Siddique, and D. Manivannan, "Effect of balancing data using synthetic data on the performance of machine learning classifiers for intrusion detection in computer networks," *IEEE Access*, vol. 10, pp. 96731–96747, 2022.
- [25] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "Smote: synthetic minority over-sampling technique," *Journal of artificial intelligence research*, vol. 16, pp. 321–357, 2002.
- [26] H. He, Y. Bai, E. A. Garcia, and S. Li, "Adasyn: Adaptive synthetic sampling approach for imbalanced learning," in *2008 IEEE international joint conference on neural networks (IEEE world congress on computational intelligence)*, pp. 1322–1328, Ieee, 2008.
- [27] E. Y. Pavlenko, I. Eremenko, and A. Fatin, "Computer network clustering methods in cybersecurity problems," *Automatic Control and Computer Sciences*, vol. 56, no. 8, pp. 957–963, 2022.
- [28] C. Chokwitthaya, Y. Zhu, S. Mukhopadhyay, and A. Jafari, "Applying the gaussian mixture model to generate large synthetic data from a small data set," in *Construction Research Congress 2020*, pp. 1251–1260, American Society of Civil Engineers Reston, VA, 2020.
- [29] L. N. Martins, F. B. Gonçalves, and T. P. Galletti, "Generation and analysis of synthetic data via bayesian networks: a robust approach for uncertainty quantification via bayesian paradigm," *arXiv preprint arXiv:2402.17915*, 2024.
- [30] D. P. Kingma, M. Welling, *et al.*, "An introduction to variational autoencoders," *Foundations and Trends® in Machine Learning*, vol. 12, no. 4, pp. 307–392, 2019.
- [31] A. Kotelnikov, D. Baranchuk, I. Rubachev, and A. Babenko, "Tabddpm: Modelling tabular data with diffusion models," in *International Conference on Machine Learning*, pp. 17564–17579, PMLR, 2023.
- [32] L. Xu, M. Skoularidou, A. Cuesta-Infante, and K. Veeramachaneni, "Modeling tabular data using conditional gan," *Advances in neural information processing systems*, vol. 32, 2019.
- [33] S. Kamthe, S. Assefa, and M. Deisenroth, "Copula flows for synthetic data generation," *arXiv preprint arXiv:2101.00598*, 2021.
- [34] Y. Zhang, N. Zaidi, J. Zhou, and G. Li, "Ganblr++: incorporating capacity to generate numeric attributes and leveraging unrestricted bayesian networks," in *Proceedings of the 2022 SIAM International Conference on Data Mining (SDM)*, pp. 298–306, SIAM, 2022.
- [35] A. Alshantti, D. Varagnolo, A. Rasheed, A. Rahmati, and F. Westad, "Castgan: Cascaded generative adversarial network for realistic tabular data synthesis," *IEEE Access*, 2024.
- [36] X. Lu, X. Ye, and Y. Cheng, "An overlapping minimization-based over-sampling algorithm for binary imbalanced classification," *Engineering Applications of Artificial Intelligence*, vol. 133, p. 108107, 2024.
- [37] O. H. Abdulganiyu, T. A. Tchakoucht, Y. K. Saheed, and H. A. Ahmed, "Xidintfl-vae: Xgboost-based intrusion detection of imbalance network traffic via class-wise focal loss variational autoencoder," *The Journal of Supercomputing*, vol. 81, no. 1, pp. 1–38, 2025.
- [38] Y. Zhang, N. Zaidi, J. Zhou, and G. Li, "Interpretable tabular data generation," *Knowledge and Information Systems*, vol. 65, no. 7, pp. 2935–2963, 2023.
- [39] A. Schoen, G. Blanc, P.-F. Gimenez, Y. Han, F. Majorczyk, and L. Me, "A tale of two methods: Unveiling the limitations of gan and the rise of bayesian networks for synthetic network traffic generation," in *2024 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*, pp. 273–286, IEEE, 2024.
- [40] D. Cullen, J. Halladay, N. Briner, R. Basnet, J. Bergen, and T. Doleck, "Evaluation of synthetic data generation techniques in the domain of anonymous traffic classification," *IEEE Access*, vol. 10, pp. 129612–129625, 2022.
- [41] F. Benali, D. Bodénès, N. Labroche, and C. de Runz, "Mtcopula: Synthetic complex data generation using copula," in *23rd International Workshop on Design, Optimization, Languages and Analytical Processing of Big Data (DOLAP)*, pp. 51–60, 2021.
- [42] Z. Li, Y. Zhao, and J. Fu, "Sync: A copula based framework for generating synthetic data from aggregated sources," in *2020 International Conference on Data Mining Workshops (ICDMW)*, pp. 571–578, IEEE, 2020.
- [43] T. J. Anande, S. Al-Saadi, and M. S. Leeson, "Generative adversarial networks for network traffic feature generation," *International Journal of Computers and Applications*, vol. 45, no. 4, pp. 297–305, 2023.
- [44] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of wasserstein gans," *Advances in neural information processing systems*, vol. 30, 2017.
- [45] Z. Lin, A. Jain, C. Wang, G. Fanti, and V. Sekar, "Using gans for sharing networked time series data: Challenges, initial promise, and open questions," in *Proceedings of the ACM Internet Measurement Conference*, pp. 464–483, 2020.
- [46] C. Pandey, V. Tiwari, A. L. Imoize, C.-T. Li, C.-C. Lee, and D. S. Roy, "5gt-gan: Enhancing data augmentation for 5g-enabled mobile edge computing in smart cities," *IEEE Access*, vol. 11, pp. 120983–120996, 2023.
- [47] N. Sivaroopan, D. Bandara, C. Madarasingha, G. Jourjon, A. P. Jayasumana, and K. Thilakarathna, "Netdiffus: Network traffic generation by diffusion models through time-series imaging," *Computer Networks*, vol. 251, p. 110616, 2024.
- [48] A. Chu, X. Jiang, S. Liu, A. Bhagoji, F. Bronzino, P. Schmitt, and N. Feamster, "Feasibility of state space models for network traffic generation," in *Proceedings of the 2024 SIGCOMM Workshop on Networks for AI Computing*, pp. 9–17, 2024.
- [49] J. Qu, X. Ma, and J. Li, "Trafficgpt: Breaking the token barrier for efficient long traffic analysis and generation," *arXiv preprint arXiv:2403.05822*, 2024.
- [50] M. Goyal and Q. H. Mahmoud, "A systematic review of synthetic data generation techniques using generative ai," *Electronics*, vol. 13, no. 17, p. 3509, 2024.
- [51] X. Guo and Y. Chen, "Generative ai for synthetic data generation: Methods, challenges and the future," *arXiv preprint arXiv:2403.04190*, 2024.
- [52] T. Liu, J. Fan, G. Li, N. Tang, and X. Du, "Tabular data synthesis with generative adversarial networks: design space and optimizations," *The VLDB Journal*, vol. 33, no. 2, pp. 255–280, 2024.
- [53] A. D'Alconzo, I. Drago, A. Morichetta, M. Mellia, and P. Casas, "A survey on big data for network traffic monitoring and analysis," *IEEE Transactions on Network and Service Management*, vol. 16, no. 3, pp. 800–813, 2019.
- [54] M. Wolf, J. Tritscher, D. Landes, A. Hotho, and D. Schlör, "Benchmarking of synthetic network data: Reviewing challenges and approaches," *Computers & Security*, p. 103993, 2024.
- [55] S. C.-H. Yang, B. Eaves, M. T. Schmidt, K. B. Swanson, and P. Shafto, "Structured evaluation of synthetic tabular data," 2023.
- [56] A. Kiran and S. S. Kumar, "A methodology and an empirical analysis to determine the most suitable synthetic data generator," *IEEE Access*, 2024.
- [57] J. Xu, L. Wei, W. Wu, A. Wang, Y. Zhang, and F. Zhou, "Privacy-preserving data integrity verification by using lightweight streaming authenticated data structures for healthcare cyber-physical system," *Future Generation Computer Systems*, vol. 108, pp. 1287–1296, 2020.
- [58] M. Tavallae, E. Bagheri, W. Lu, and A. A. Ghorbani, "A detailed analysis of the kdd cup 99 data set," in *2009 IEEE symposium on computational intelligence for security and defense applications*, pp. 1–6, Ieee, 2009.
- [59] A. F. Karr, C. N. Kohnen, A. Oganian, J. P. Reiter, and A. P. Sanil, "A framework for evaluating the utility of data altered to protect confidentiality," *The American Statistician*, vol. 60, no. 3, pp. 224–232, 2006.
- [60] M. Pereira, M. Kshirsagar, S. Mukherjee, R. Dodhia, J. Lavista Ferres, and R. de Sousa, "Assessment of differentially private synthetic data for utility and fairness in end-to-end machine learning pipelines for tabular data," *Plos one*, vol. 19, no. 2, p. e0297271, 2024.

- [61] J. Snoke, G. M. Raab, B. Nowok, C. Dibben, and A. Slavkovic, "General and specific utility measures for synthetic data," *Journal of the Royal Statistical Society Series A: Statistics in Society*, vol. 181, no. 3, pp. 663–688, 2018.
- [62] P. Thölke, Y.-J. Mantilla-Ramos, H. Abdelhedi, C. Maschke, A. Dehgan, Y. Harel, A. Kemtur, L. M. Berrada, M. Sahraoui, T. Young, *et al.*, "Class imbalance should not throw you off balance: Choosing the right classifiers and performance metrics for brain decoding with imbalanced data," *NeuroImage*, vol. 277, p. 120253, 2023.
- [63] E. R. Fernandes and A. C. de Carvalho, "Evolutionary inversion of class distribution in overlapping areas for multi-class imbalanced learning," *Information Sciences*, vol. 494, pp. 141–154, 2019.
- [64] C. E. Shannon, "A mathematical theory of communication," *The Bell system technical journal*, vol. 27, no. 3, pp. 379–423, 1948.
- [65] T. M. Cover, *Elements of information theory*. John Wiley & Sons, 1999.
- [66] W. McGill, "Multivariate information transmission," *Transactions of the IRE Professional Group on Information Theory*, vol. 4, no. 4, pp. 93–111, 1954.
- [67] S. Manzoor, M. K. Siddiqui, and S. Ahmad, "On entropy measures of molecular graphs using topological indices," *Arabian Journal of Chemistry*, vol. 13, no. 8, pp. 6285–6298, 2020.
- [68] K. Church and P. Hanks, "Word association norms, mutual information, and lexicography," *Computational linguistics*, vol. 16, no. 1, pp. 22–29, 1990.
- [69] I. Guyon and A. Elisseeff, "An introduction to variable and feature selection," *Journal of machine learning research*, vol. 3, no. Mar, pp. 1157–1182, 2003.
- [70] L. Breiman, "Random forests," *Machine learning*, vol. 45, pp. 5–32, 2001.
- [71] D. Stiawan, M. Y. B. Idris, A. M. Bamhdi, R. Budiarto, *et al.*, "Cicids-2017 dataset feature analysis with information gain for anomaly detection," *IEEE Access*, vol. 8, pp. 132911–132921, 2020.
- [72] R. Mahto, S. U. Ahmed, R. u. Rahman, R. M. Aziz, P. Roy, S. Mallik, A. Li, and M. A. Shah, "A novel and innovative cancer classification framework through a consecutive utilization of hybrid feature selection," *BMC bioinformatics*, vol. 24, no. 1, p. 479, 2023.
- [73] X. Zhang, X.-M. Zhao, K. He, L. Lu, Y. Cao, J. Liu, J.-K. Hao, Z.-P. Liu, and L. Chen, "Inferring gene regulatory networks from gene expression data by path consistency algorithm based on conditional mutual information," *Bioinformatics*, vol. 28, no. 1, pp. 98–104, 2012.
- [74] J. Hua, W. D. Tembe, and E. R. Dougherty, "Performance of feature-selection methods in the classification of high-dimension data," *Pattern Recognition*, vol. 42, no. 3, pp. 409–424, 2009.
- [75] H. Peng, F. Long, and C. Ding, "Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 27, no. 8, pp. 1226–1238, 2005.
- [76] S. D. Team, "Sdv user guides: Single table - copulagan." [https://sdv.dev/SDV/user\\_guides/single\\_table/copulagan.html](https://sdv.dev/SDV/user_guides/single_table/copulagan.html), 2022.
- [77] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.
- [78] D. A. Reynolds *et al.*, "Gaussian mixture models.," *Encyclopedia of biometrics*, vol. 741, no. 659-663, 2009.

## APPENDIX

In this appendix, we present the complete set of probability distribution (PD) comparison curves for the NSL-KDD (26 selected pre-processed variables) and CIC-IDS2017 (21 selected pre-processed variables) datasets. These graphs showcase the PDs for each variable across twelve different sampling methods mentioned in the section III.

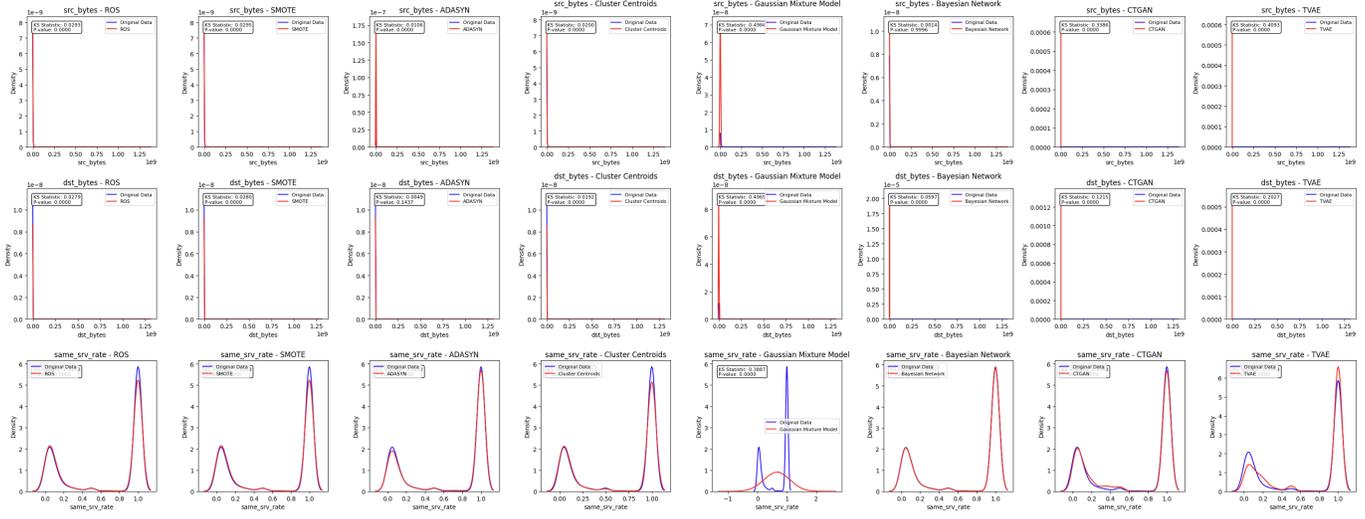


Fig. 5. Probability Distribution Comparison for NSL KDD (Part 1/18)

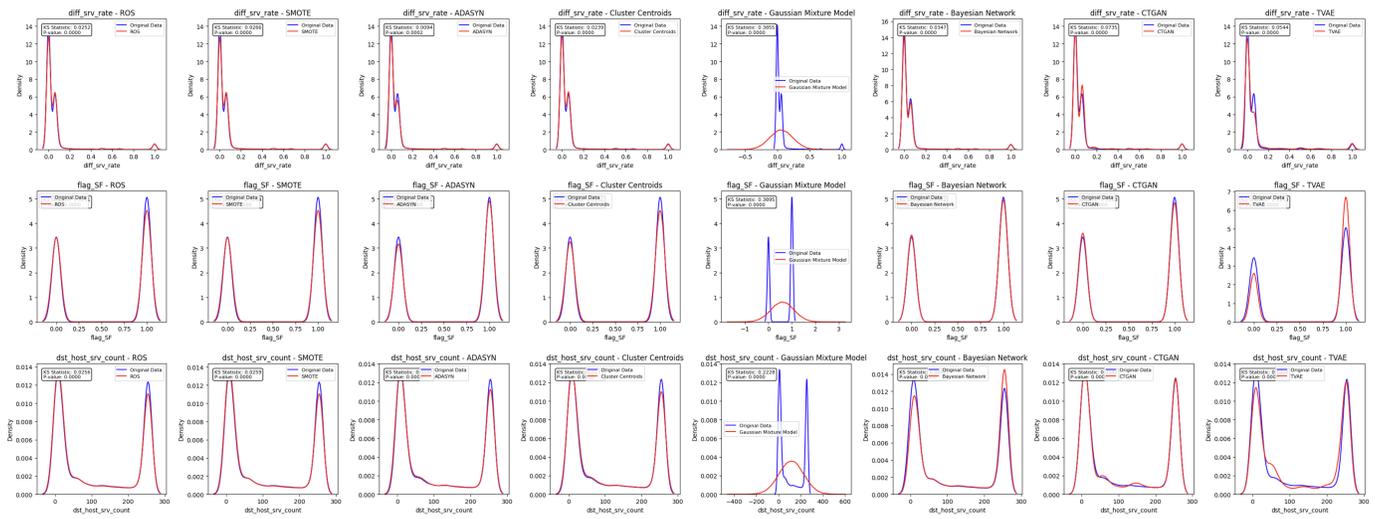


Fig. 6. Probability Distribution Comparison for NSL KDD (Part 2/18)

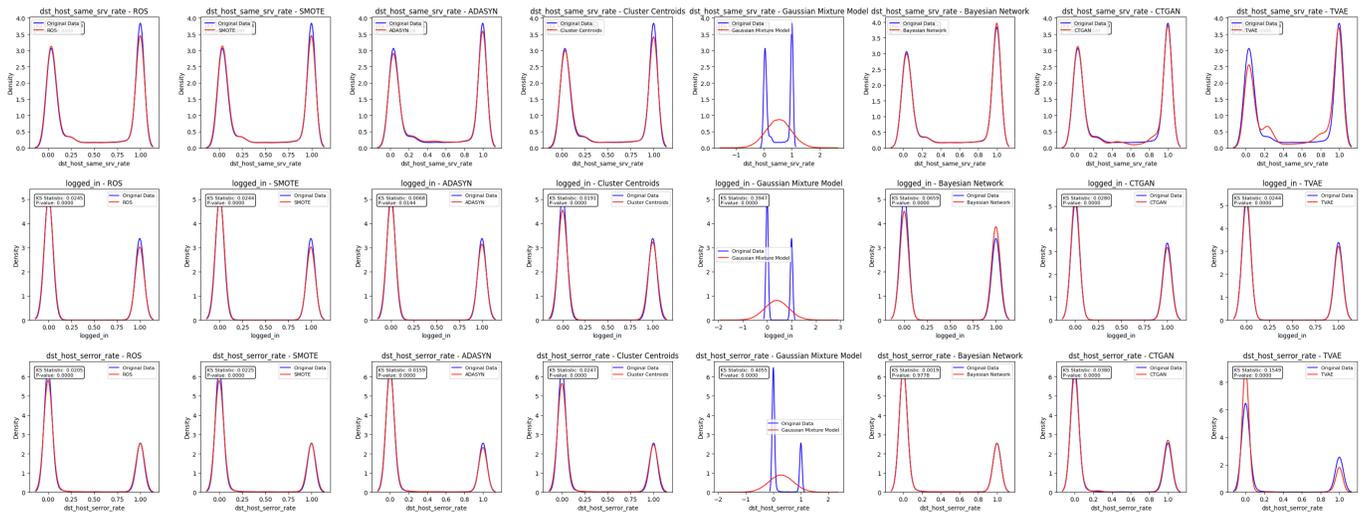


Fig. 7. Probability Distribution Comparison for NSL KDD (Part 3/18)

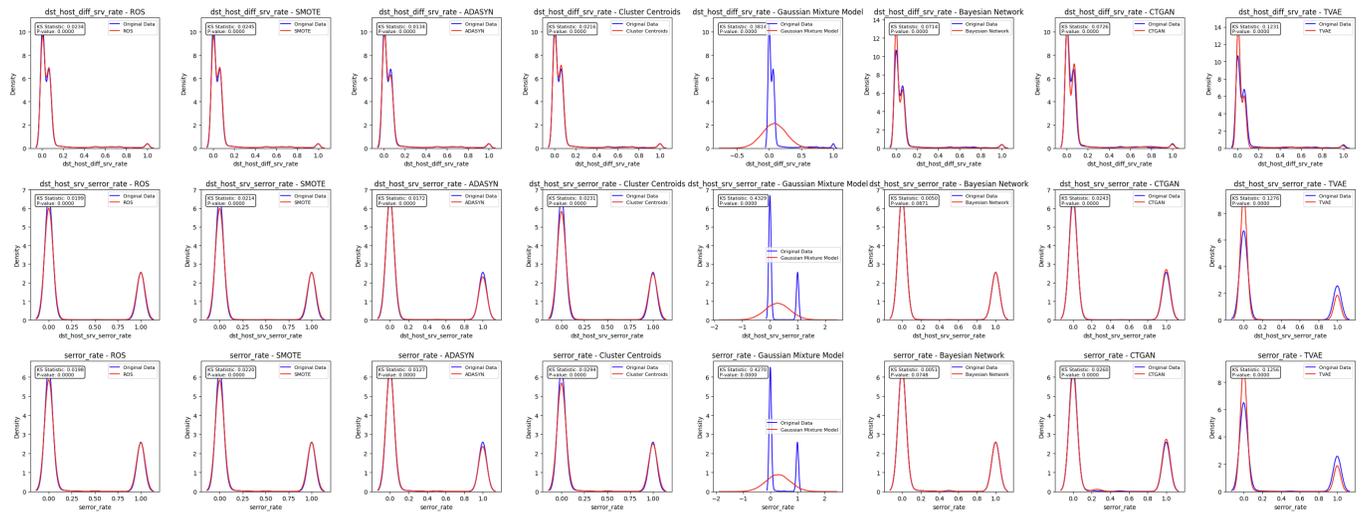


Fig. 8. Probability Distribution Comparison for NSL KDD (Part 4/18)

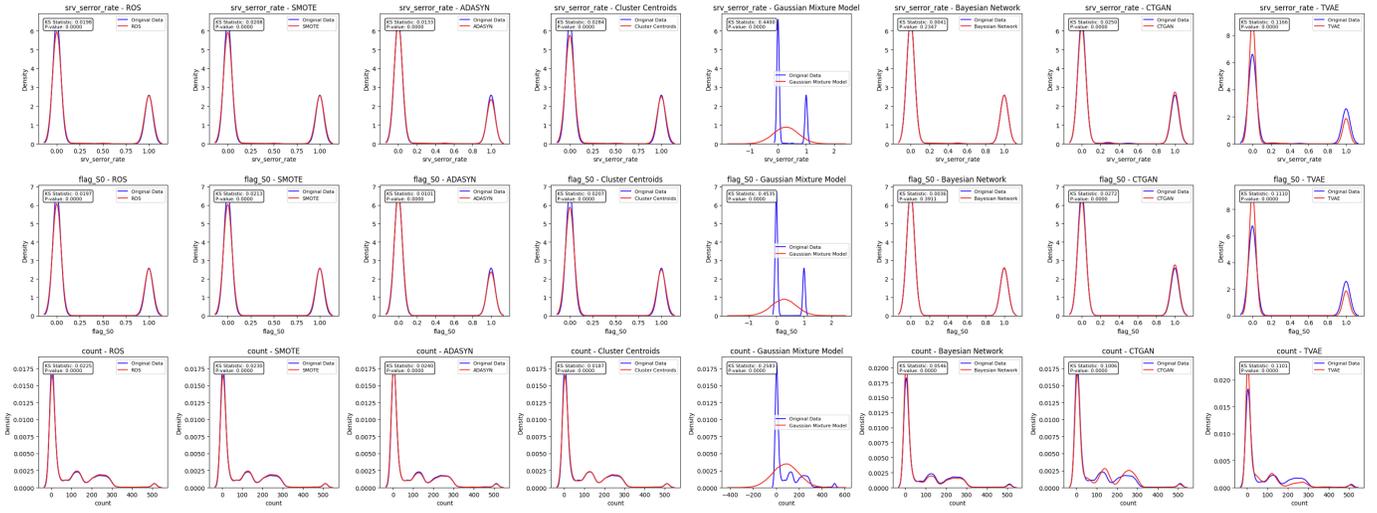


Fig. 9. Probability Distribution Comparison for NSL KDD (Part 5/18)

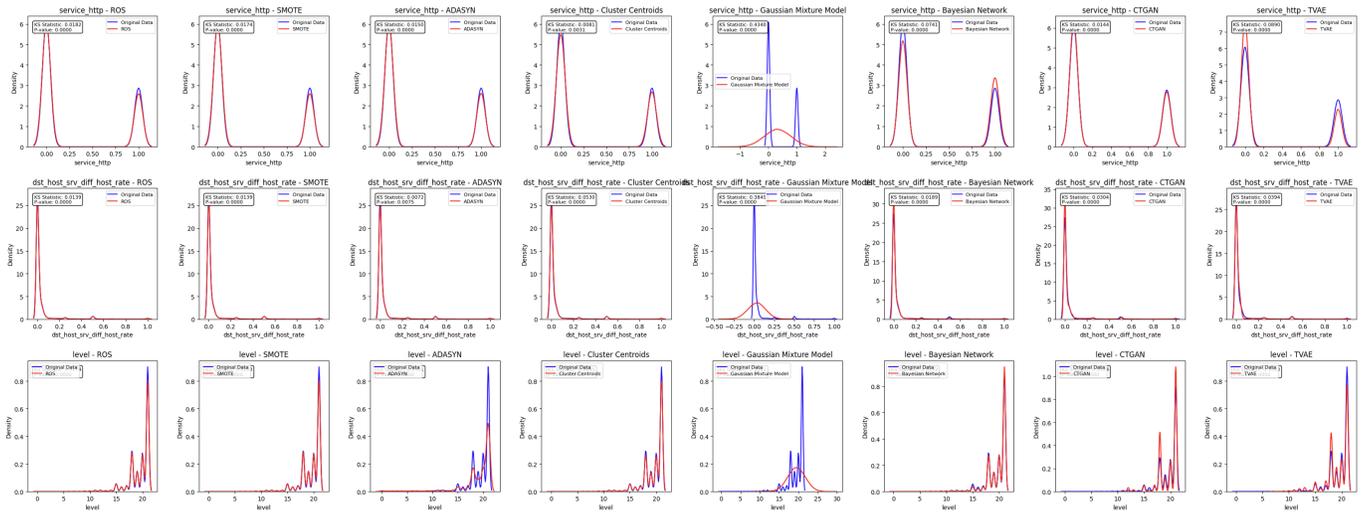


Fig. 10. Probability Distribution Comparison for NSL KDD (Part 6/18)

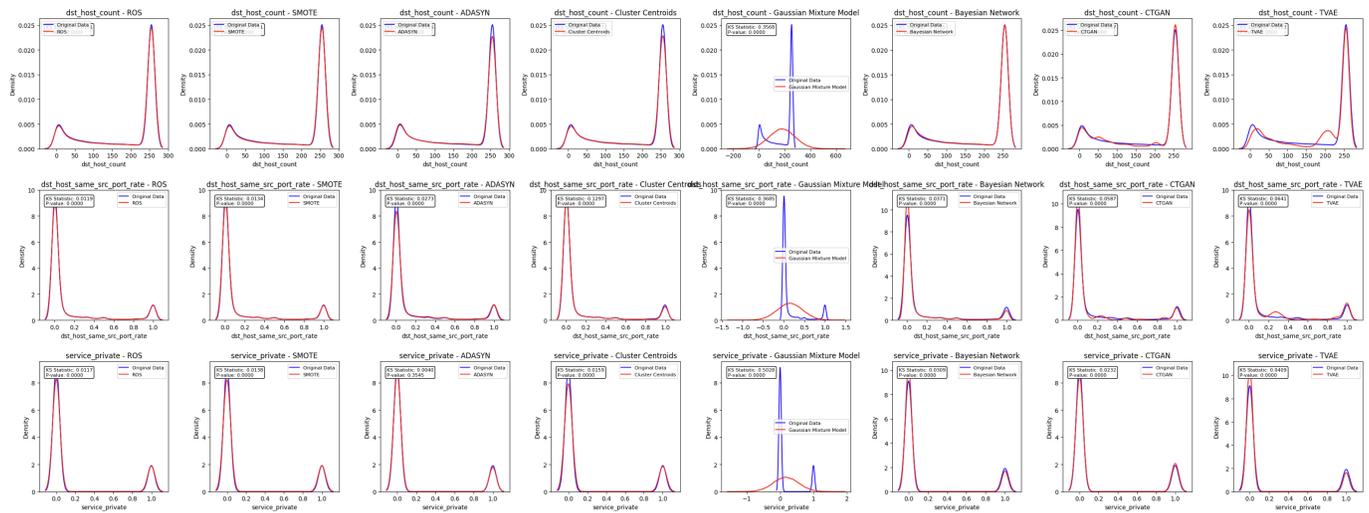


Fig. 11. Probability Distribution Comparison for NSL KDD (Part 7/18)

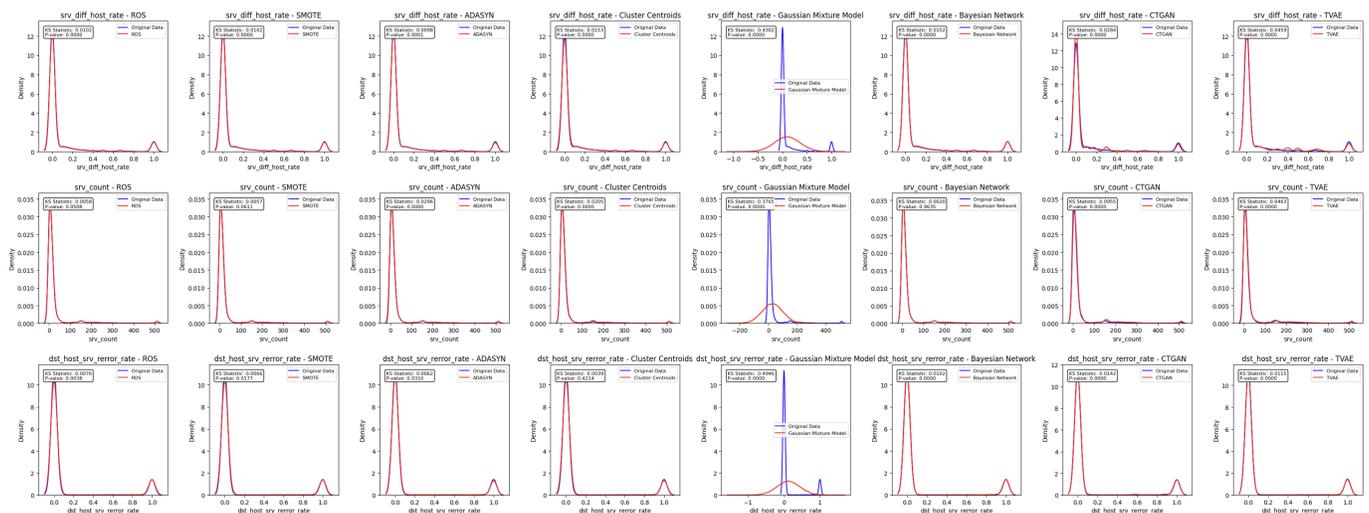


Fig. 12. Probability Distribution Comparison for NSL KDD (Part 8/18)

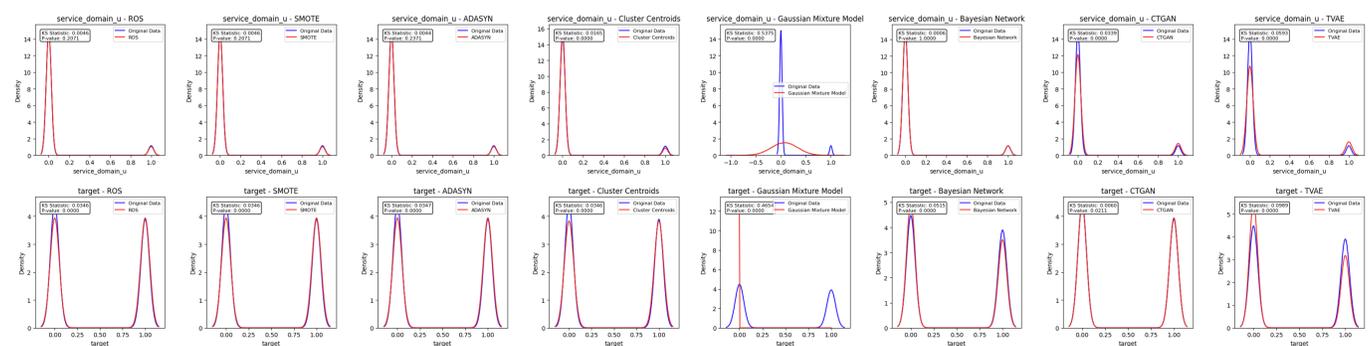


Fig. 13. Probability Distribution Comparison for NSL KDD (Part 9/18)

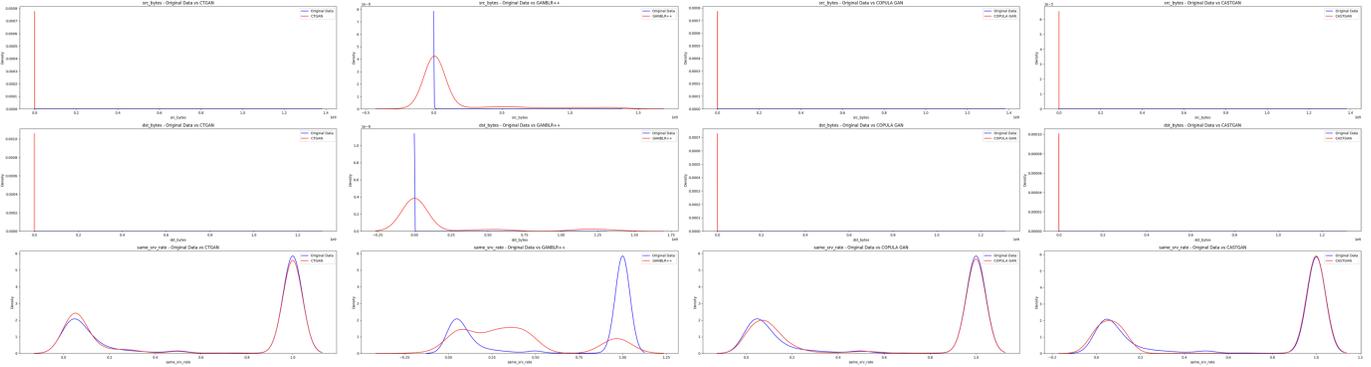


Fig. 14. Probability Distribution Comparison for NSL KDD (Part 10/18)

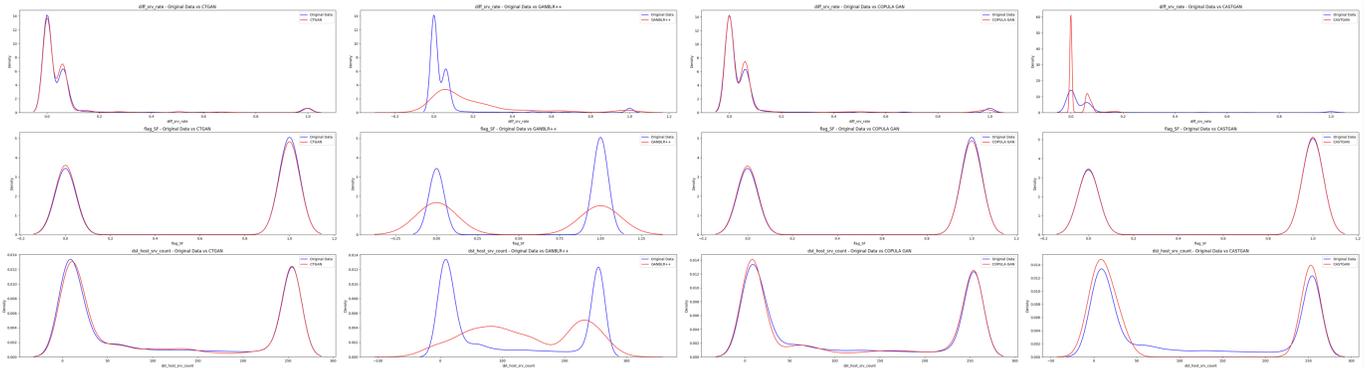


Fig. 15. Probability Distribution Comparison for NSL KDD (Part 11/18)

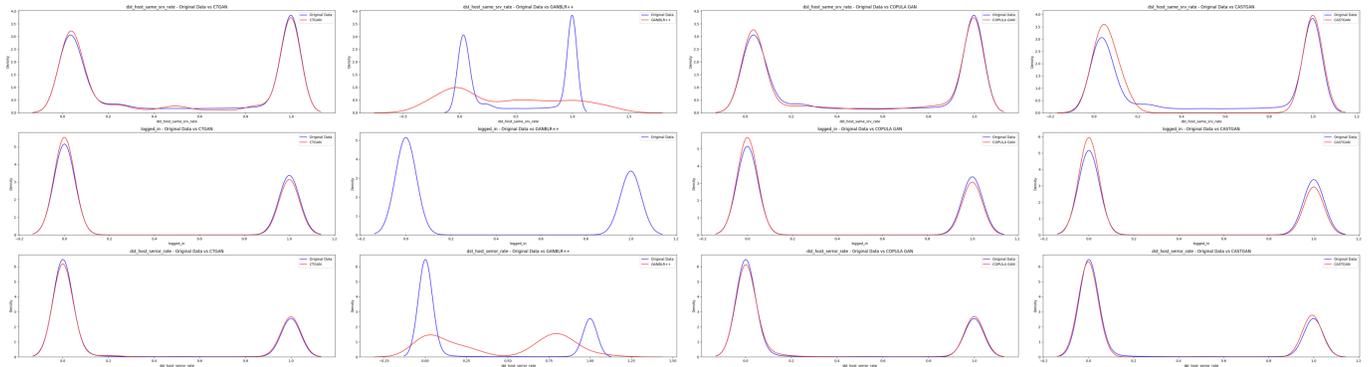


Fig. 16. Probability Distribution Comparison for NSL KDD (Part 12/18)

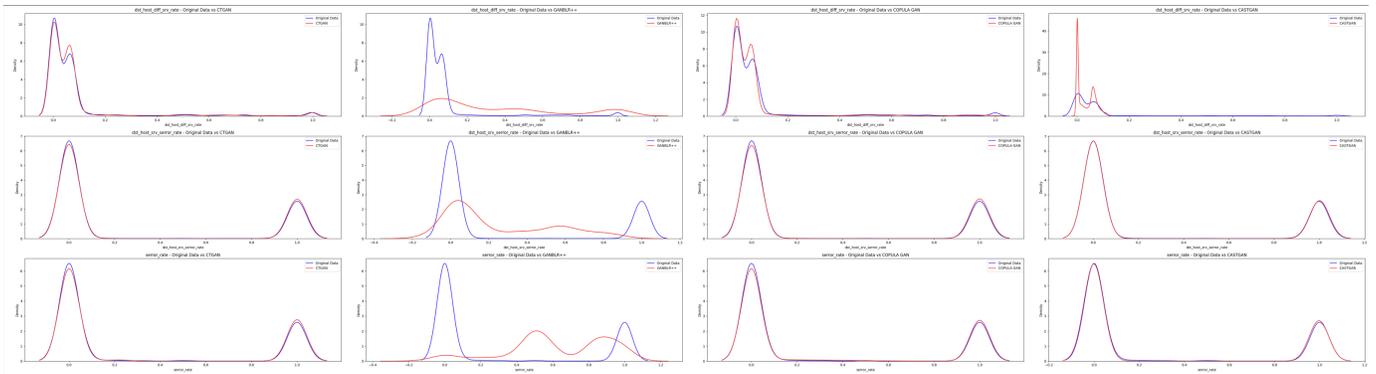


Fig. 17. Probability Distribution Comparison for NSL KDD (Part 13/18)

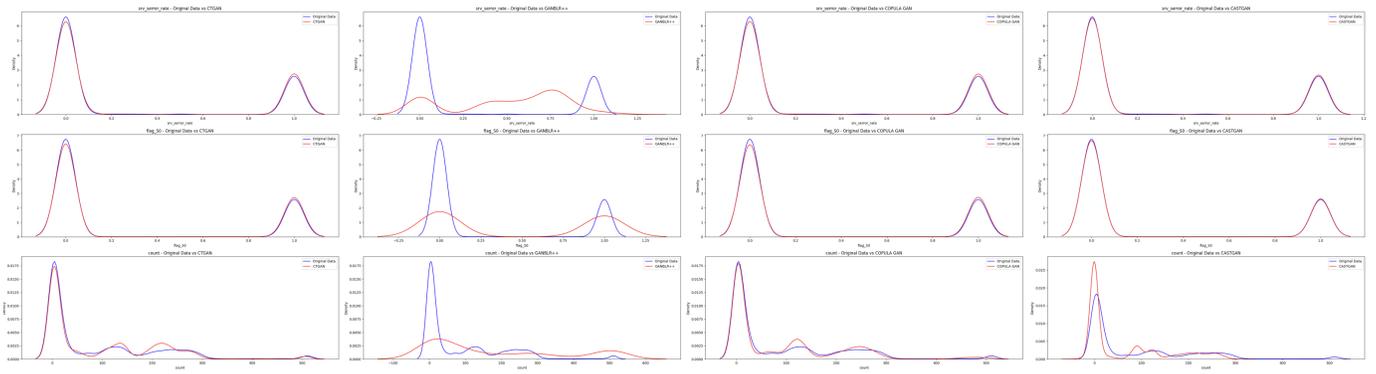


Fig. 18. Probability Distribution Comparison for NSL KDD (Part 14/18)

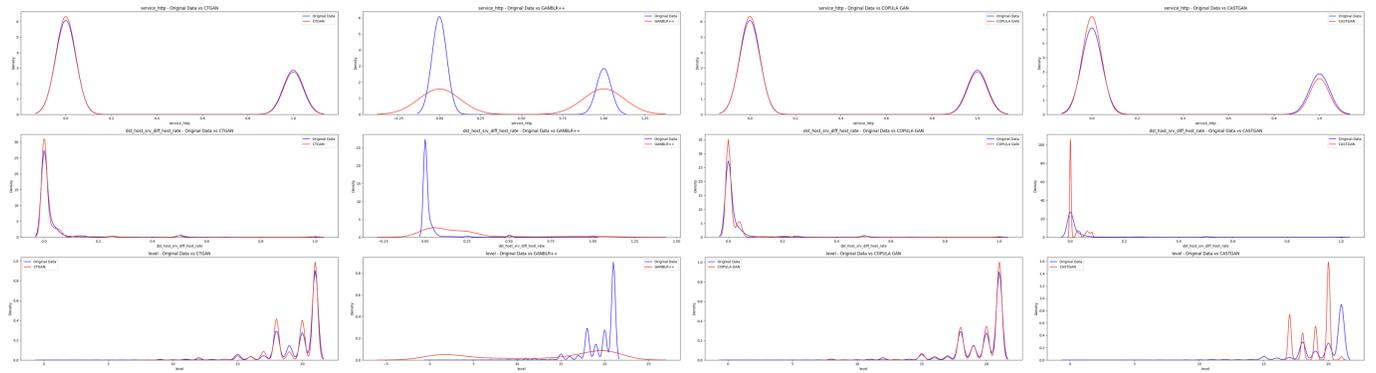


Fig. 19. Probability Distribution Comparison for NSL KDD (Part 15/18)

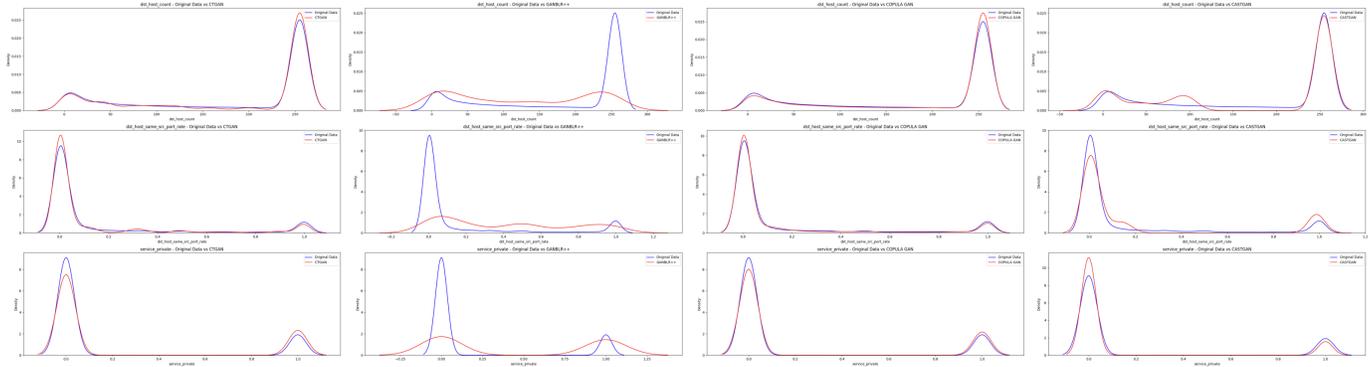


Fig. 20. Probability Distribution Comparison for NSL KDD (Part 16/18)

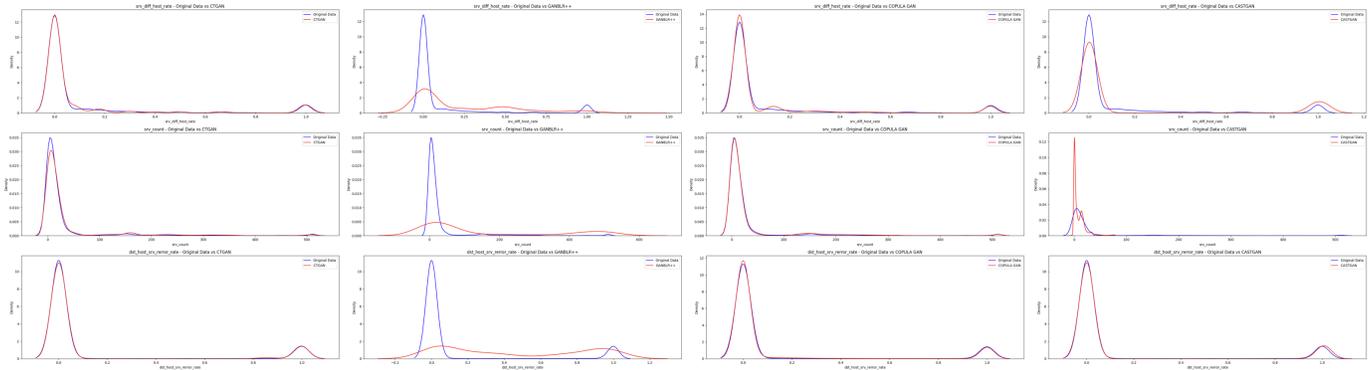


Fig. 21. Probability Distribution Comparison for NSL KDD (Part 17/18)

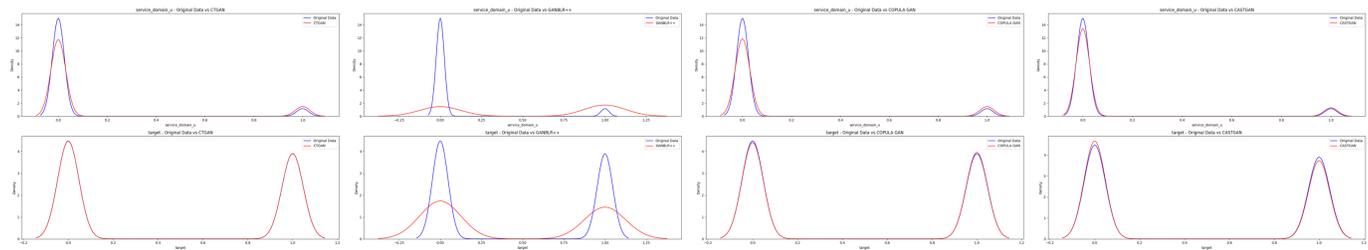


Fig. 22. Probability Distribution Comparison for NSL KDD (Part 18/18)

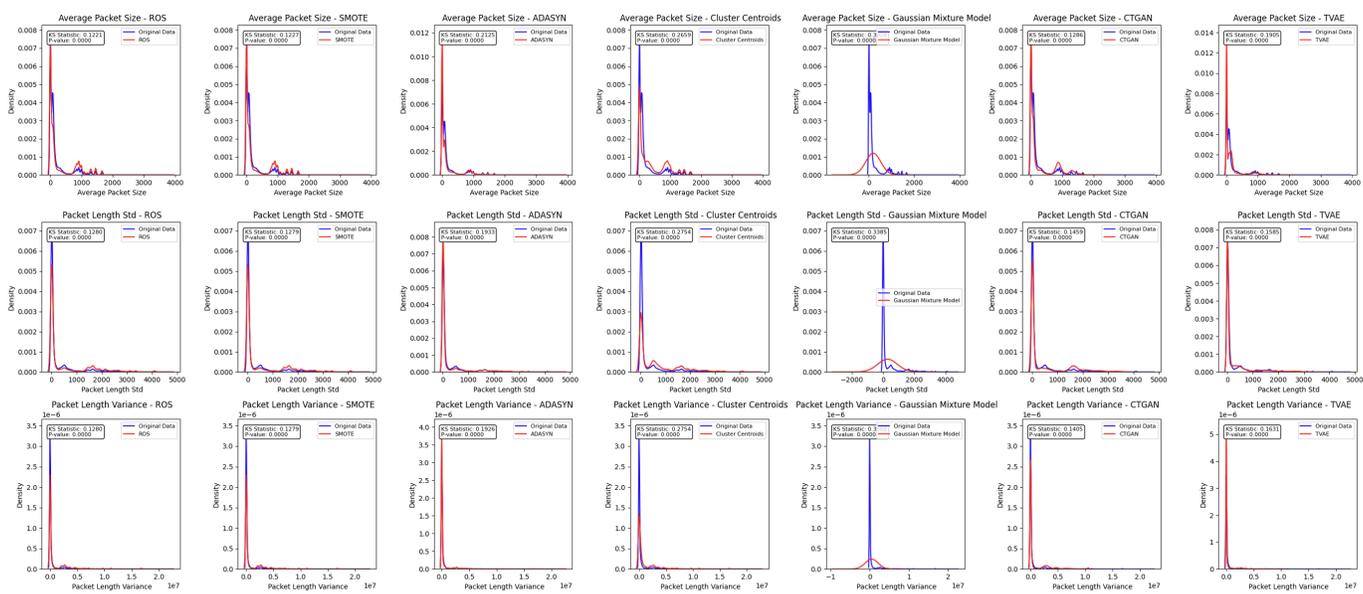


Fig. 23. Probability Distribution Comparison for CIC-IDS2017 (Part 1/7)

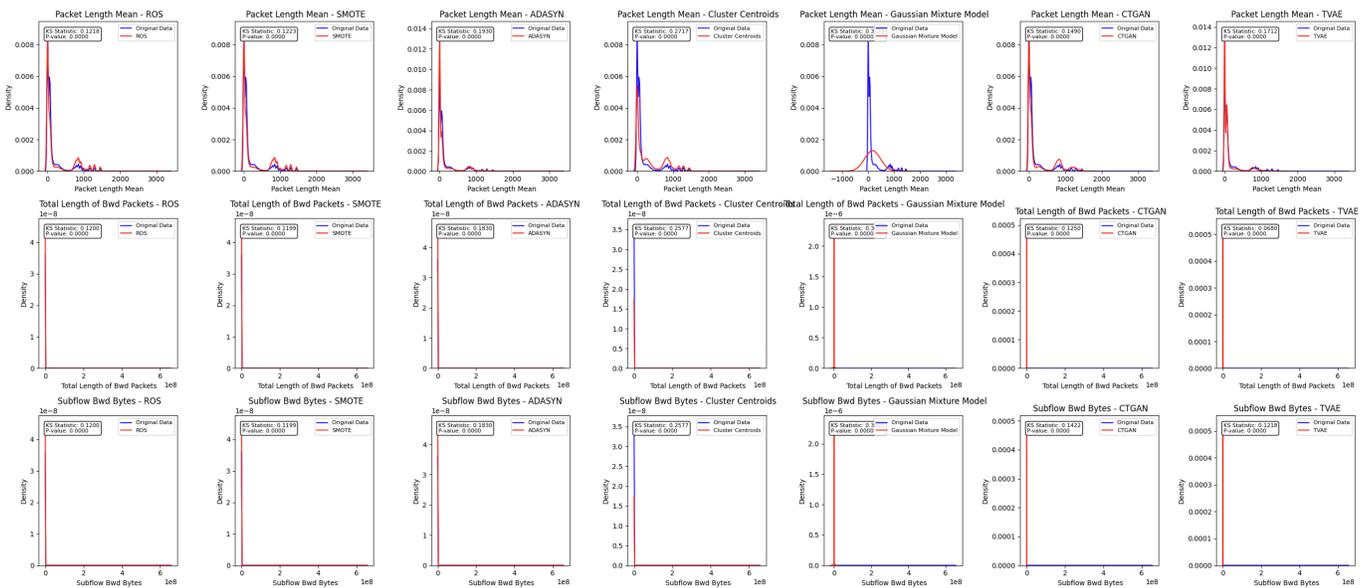


Fig. 24. Probability Distribution Comparison for CIC-IDS2017 (Part 2/7)

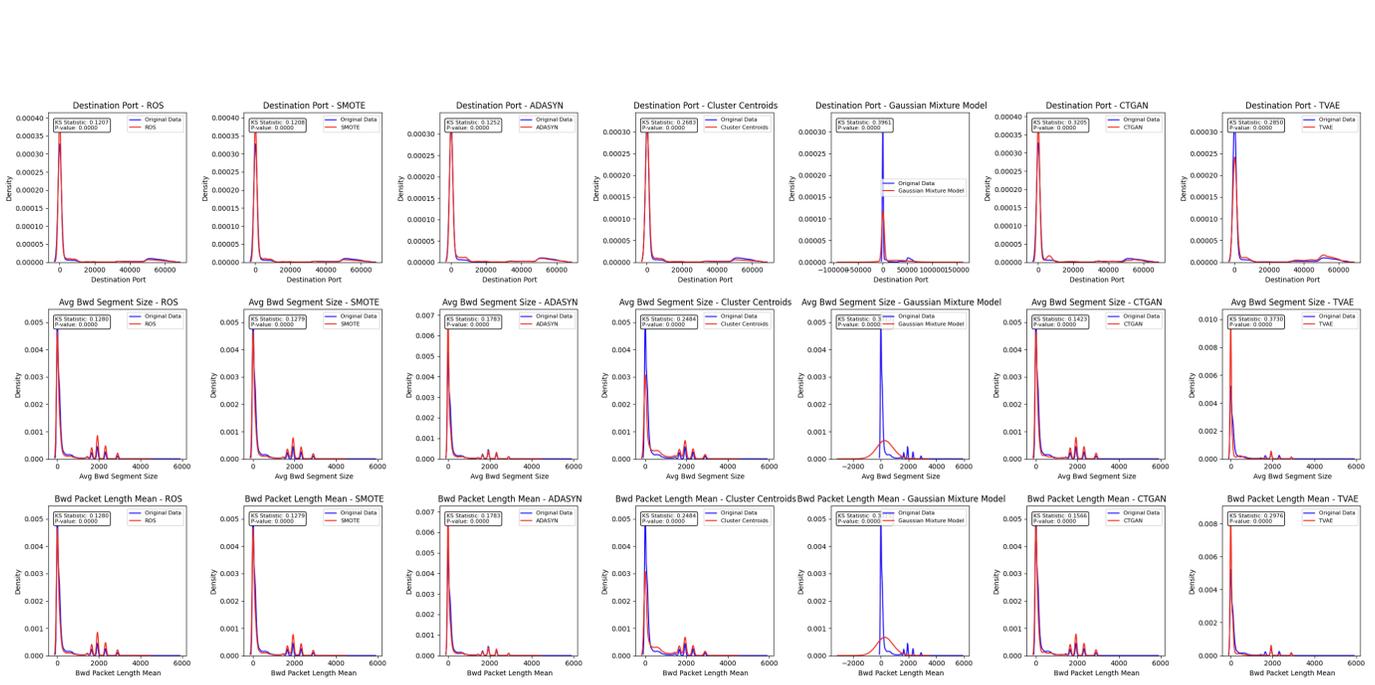


Fig. 25. Probability Distribution Comparison for CIC-IDS2017 (Part 3/7)

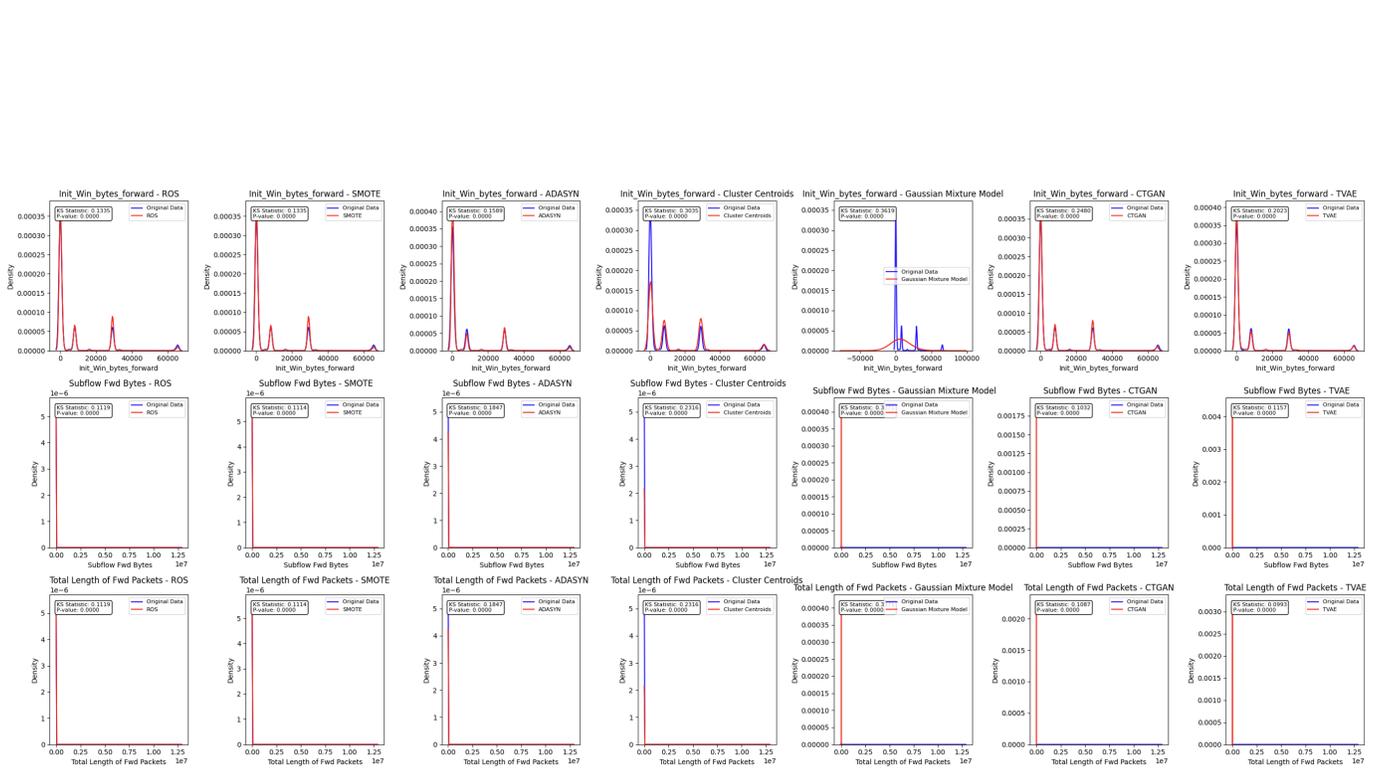


Fig. 26. Probability Distribution Comparison for CIC-IDS2017 (Part 4/7)

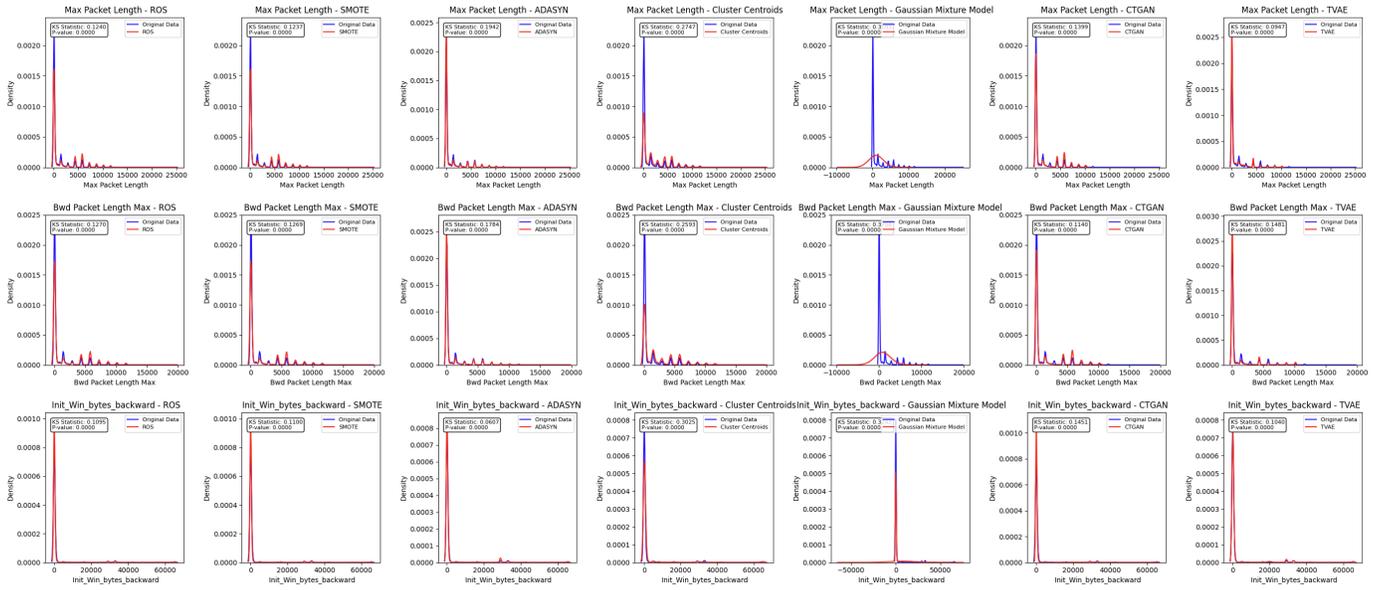


Fig. 27. Probability Distribution Comparison for CIC-IDS2017 (Part 5/7)

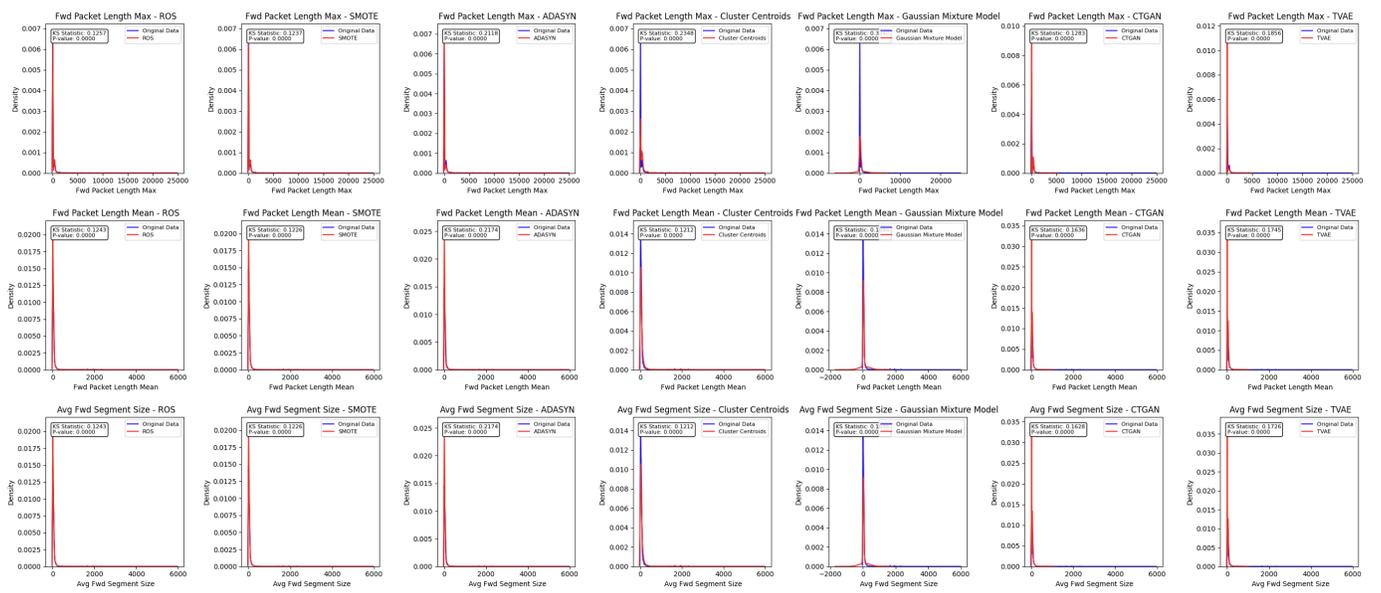


Fig. 28. Probability Distribution Comparison for CIC-IDS2017 (Part 6/7)

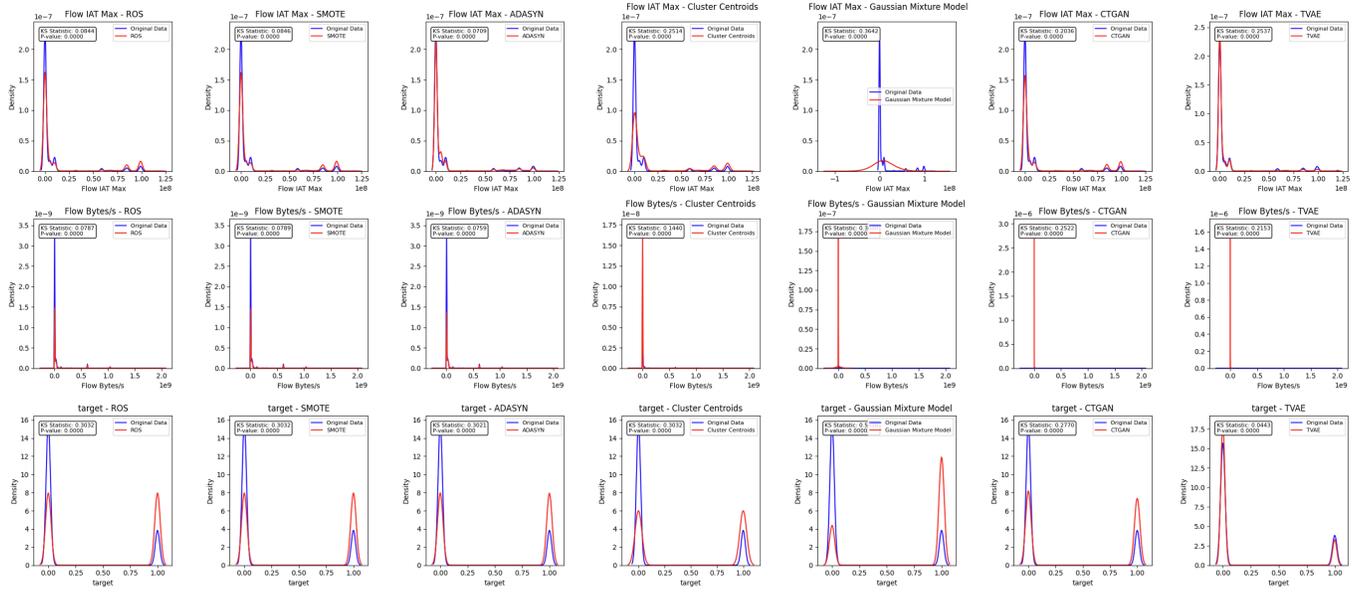


Fig. 29. Probability Distribution Comparison for CIC-IDS2017 (Part 7/7)

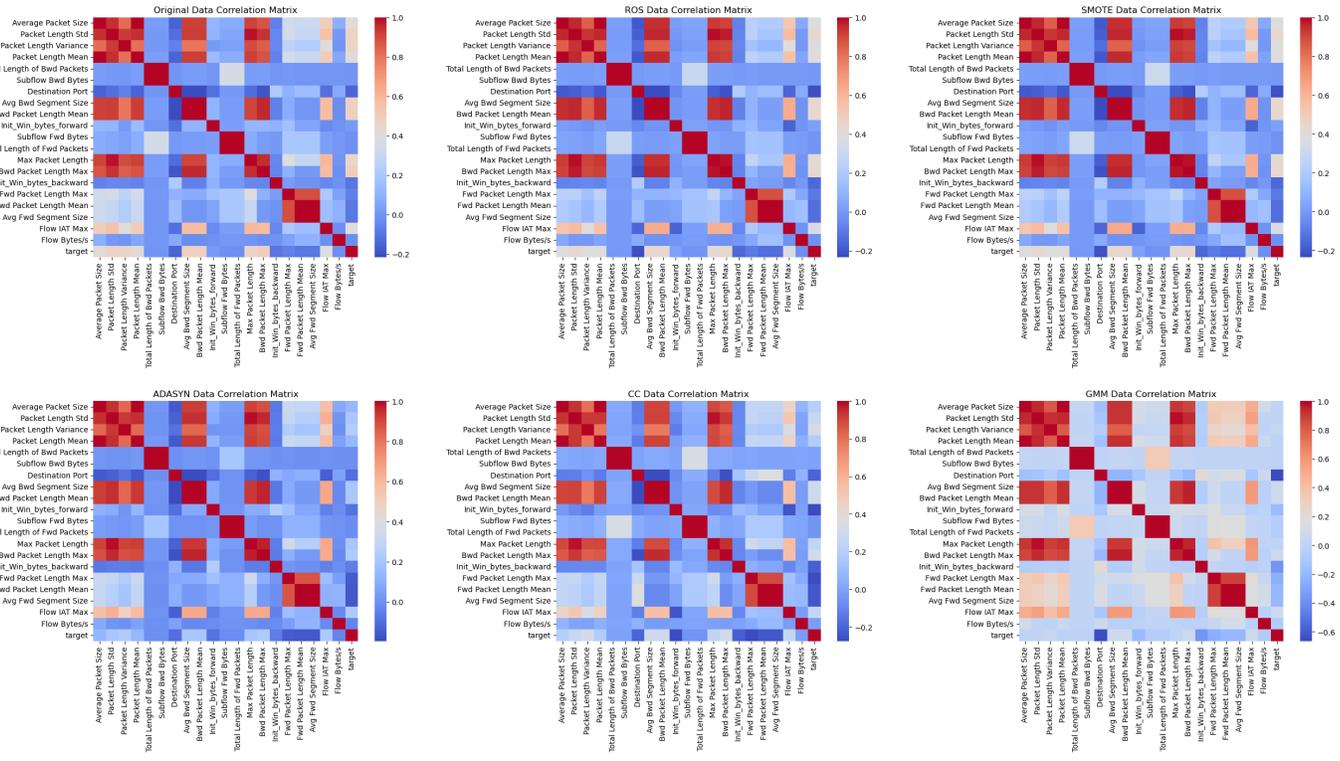


Fig. 30. Correlation Heatmap for CIC-IDS2017 (Part 1/2)

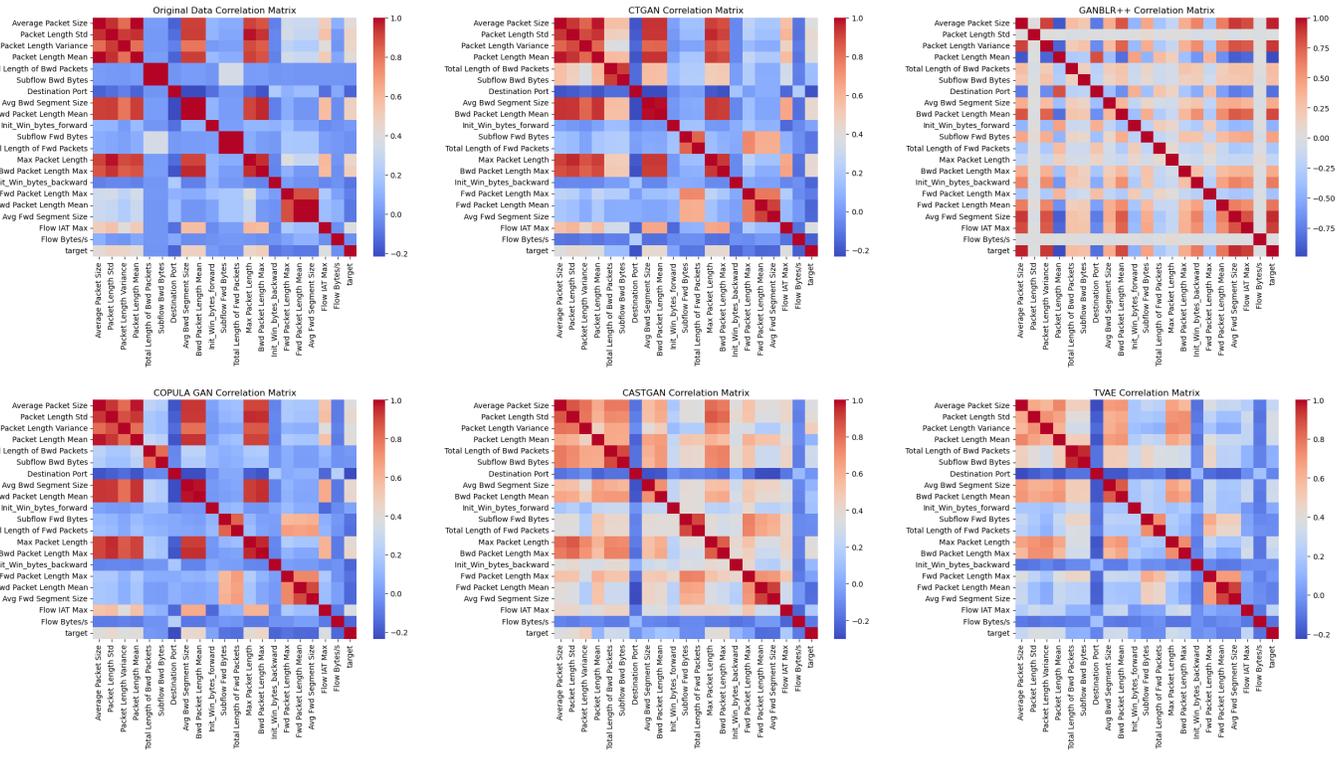


Fig. 31. Correlation Heatmap for CIC-IDS2017 (Part 2/2)