# Applying Generative mock Neuro Forge Networks for Synthetic Data Generation in AI Healthcare Systems

Md Mahadi Hasan<sup>1</sup>, Seaam Bin Masud<sup>2</sup>, Md Rafiuddin Siddiky<sup>3</sup>, Samia Ara Chowdhury<sup>4</sup>, Intiser Islam<sup>5</sup>, Israt Jahan<sup>6</sup>,

1.M.S. in Information Technology, Washington University of Science and Technology, Alexandria, Virginia, United States, mahadih 384@qmail.com

ORCID: https://orcid.org/0009-0003-1010-7442

2.M.S. in Information Technology Project Management, Wilmington University, New Castle, Delaware, United States, seaam.masud@gmail.com

ORCID: https://orcid.org/0009-0001-3725-4042

3.M.S. in Information Systems Technology, Wilmington University, New Castle, Delaware, United States. mdrafiuddinsiddiky@gmail.com

ORCID: https://orcid.org/0009-0006-6528-6486

4.M.S. in Information Technology, ST FRANCIS COLLEGE , Brooklyn, New York, USA, Samia.r.chowdhury@gmail.com

5.M.S. in Computer Science, School of Engineering, University of Bridgeport, Bridgeport, Connecticut, USA, iislam@my.bridgeport.edu

6. Washington University of Science and Technology, Master of Science in Information Technology, Virginia, USA, ijahan.student@wust.edu

### Abstract-

The rapid advancement of artificial intelligence (AI) in healthcare has catalyzed the need for large, diverse, and high-quality datasets to train robust machine learning models. However, acquiring real-world medical data presents challenges due to privacy concerns, regulatory restrictions, and data scarcity. Generative mock NeuroForge Networks (GMNFNs) offer a promising solution by enabling synthetic data generation that mimics the complexity and variability of real-world datasets while preserving patient confidentiality. This paper introduces a novel three-step framework for synthetic data generation in AI healthcare systems: (1) HoloScope Sampling—a pre-processing algorithm that ensures input data diversity and represents the full spectrum of real-world scenarios; (2) Generative mock NeuroForge Networks (GMNFNs)—a cutting-edge architecture designed to generate high-fidelity synthetic datasets while addressing privacy and ethical constraints; and (3) Fuzzy press DataTrust Validator (FPDTV)—a postgeneration algorithm that quantitatively evaluates the reliability and utility of synthetic datasets using advanced statistical and domain-specific metrics. By integrating these steps, this research demonstrates a pathway to bridge data gaps, enhance model performance, and mitigate biases in healthcare AI systems. Ethical considerations and the integration of these algorithms into existing frameworks are discussed, providing a roadmap for accelerating innovation while adhering to privacy and regulatory standards.

Index Terms—Artificial Intelligence, Healthcare, Generative mock NeuroForge Networks, Synthetic data generation

# I. INTRODUCTION

Artificial intelligence (AI) has emerged as a powerful tool for addressing critical challenges in healthcare, revolutionizing areas such as disease diagnosis, treatment planning, drug discovery, and patient monitoring. AI systems rely heavily on high-quality, diverse, and extensive datasets to achieve accuracy, reliability, and generalizability. However, the healthcare sector faces unique challenges in acquiring such datasets. Privacy concerns, ethical considerations, and strict regulatory frameworks, such as GDPR and HIPAA, limit access to sensitive patient data. Moreover, data scarcity and inherent biases in available datasets further hinder the development of equitable and effective AI solutions.

Synthetic data generation has become a promising approach to overcome these barriers. By creating artificial datasets that replicate the statistical properties and complexity of real-world healthcare data, synthetic data enables researchers and developers to train AI models without exposing sensitive information. This approach also facilitates the creation of datasets that address underrepresented populations, ensuring a more equitable and inclusive foundation for healthcare AI systems. While existing generative models, such as Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs), have made significant strides in

synthetic data generation, they face limitations in capturing domain-specific nuances and addressing the stringent ethical and regulatory requirements of the healthcare industry.

To address these challenges, this research introduces a novel three-step framework for synthetic data generation tailored specifically for healthcare applications. The first step, **HoloScope Sampling**, focuses on creating a diverse and representative input dataset by leveraging advanced sampling techniques. This ensures the inclusion of rare medical conditions, demographic diversity, and edge cases critical for building robust AI systems. The second step employs Generative mock **NeuroForge Networks**, a state-of-the-art architecture designed to generate high-fidelity synthetic healthcare datasets. This novel model incorporates domain-specific knowledge and privacy-preserving mechanisms to balance data fidelity with compliance and ethical considerations. The final step, fuzzy press **DataTrust Validator (FPDTV)**, evaluates the synthetic datasets using advanced statistical, domain-specific, and ethical metrics to ensure their reliability, usability, and adherence to healthcare standards.

This framework not only addresses the pressing challenges of data scarcity and privacy in healthcare AI but also sets a benchmark for generating synthetic data that aligns with real-world applications. By bridging gaps in data availability and reducing biases, the proposed methods enable the development of AI models that are more inclusive, accurate, and impactful. The integration of this framework into healthcare workflows has the potential to accelerate AI innovation, democratize access to healthcare solutions, and improve patient outcomes globally.

The subsequent sections of this paper provide a detailed exploration of the proposed framework, its algorithms, and its application to real-world healthcare scenarios. This work contributes to the ongoing efforts to make healthcare AI systems more ethical, efficient, and effective.

The paper is organized into several sections. The **Introduction** outlines the background, challenges, and motivation for synthetic data generation in healthcare AI and introduces the proposed framework. The **Related Work** section reviews existing methodologies, highlighting limitations in privacy, data quality, and domain-specific adaptability. The **Proposed Framework** describes the **Methodology** that details the implementation and integration of these components, while the **Results and Discussion** present experimental findings, performance comparisons, and implications for healthcare AI. The **Conclusion and Future Work** summarize contributions and propose directions for further research.

# II. RELATED WORKS

As shown in earlier research [1]-[7], [16], [7]-[17], the use of synthetic datasets using learning algorithms has been investigated extensively using different data generators. Synthetic data has also been shown to be beneficial in machine learning [5, 6, 18]. In addition, other pioneers have developed data generating methods using GAN since Ian Goodfellow's creation of Generative Adversarial Networks [19] [20], [21], [22], [23].

The generating and discriminative models are the two pillars of a GAN. The first one tries to make the original data seem the same by adding noise, while the second one checks how close the created data is to the original data by comparing their distributions. Problems arise in the classification model due to GAN's inclusion of the class label as an extra attribute. Fortunately, this problem is addressed by the Conditional GAN (CGAN) [21], which improves the data quality for the classification model by treating the class labels individually. Because GAN makes use of deep networks, the system tries to retrieve the training data and enhance it till it becomes as near as possible to the original. Because of this, any private or sensitive information in the original data might be exposed. This problem was solved by DPGAN [24], which integrated GAN with differential privacy (DP). To create differentially private synthetic data, the DPGAN trains the discriminator using data that has been artificially infected with noise and then uses this data to inform its predictions. By incorporating the Private Aggregation of Teacher Ensembles (PATE) into the GAN, the PATE-GAN [20] paradigm

expands upon this concept. In order to train the discriminator, PATE-GAN uses a teacher and student model to generate a noisy dataset. This approach outperforms DPGAN when it comes to disclosure control. Differential Privacy (DP) is a tried and true method for protecting the confidentiality of publicly available datasets. Nevertheless, it impacts the usefulness and equity of artificial datasets as well. Because of this, studying how DP affects synthetic data is essential, especially in cases where there are gaps between the dataset's majority and minority classes. The uneven influence of DP on balancing classes in synthetic datasets has been shown using several kinds of Generative Adversarial Networks (GANs). To demonstrate the distinctions, Ganev et al. used PrivBayes, W-GAN, and PATE-GAN to create synthetic data from three separate datasets. Based on their results, they concluded that PATE-GAN widens the gap between minority and majority classes, whereas PrivBayes narrows it. Nevertheless, W-GAN yielded contradictory findings. Synthetic data creation for healthcare datasets utilising GAN employing DP yielded comparable findings and observations [15].

In order to gauge the integrity of artificial healthcare data, Karan Bhanot et al. dug further into formulating a strong measure. During data creation, they also insisted that fairness must be included into the dataset. The unfavourable results of the machine learning models trained on these datasets are believed to be caused by additive noise, gradient clipping, and DP Stochastic Gradient Descent, all of which introduce bias into the dataset and thereby affect fairness [11], [12]. A bigger imbalance in the synthetic data set may be produced by DP, even if the training data set includes a minor difference across the classes [25]. For more accurate results, it's best to pre-process the dataset using multi-label under-sampling until the minority and majority numbers are equal [26]. The impact of pre-processing with four distinct SDGs was shown by Blake Bullwinkel et al. The data generators employed were MST, DP-CTGAN, PATECTGAN, and SN-synth, which is a component of the Smartnoise-synth package. The research found that generators based on GANs provide different outcomes. Hyperparameter tweaking was proposed by Blake Bullwinkel et al. When the privacy budget e was equal to or more than 3.0, PATE-GAN produced better outcomes, whereas DPCTGAN produced better results when e was less than or equal to 1.0, according to the attempts to benchmark four distinct differentially private GAN-based SDGs [23]. Neither of the generators is deemed superior in the conclusion. Options for reducing bias using the privatised likelihood ratio show that the SDGs still have problems with bias [27]. In [28] the author compares and contrasts five ML algorithms—Logistic Regression, Support Vector Machines (SVM), Random Forest, Naive Bayes, and Gradient Boosting—used for sentiment categorisation in social media material. We tested these algorithms' ability to identify positive, negative, and neutral attitudes in a dataset consisting of one hundred thousand tweets gathered over a span of three months. Cleaning, normalising, and resolving class imbalance using SMOTE were among the several preprocessing steps performed on the data. We found that Logistic Regression and SVM performed equally well across all sentiment classifications, with an overall accuracy of 86.22%. With an accuracy of 82.59%, Random Forest was right behind, while Gradient Boosting and Naive Bayes also performed well, although at lesser levels (69.96% and 70.45%, respectively). In [29] the author investigates the potential of combining the distributed ledger technology of blockchain with the predictive analytics of ML to safeguard monetary transactions. Academic databases such as IEEE Xplore, Google Scholar, Scopus, Web of Science, DOAJ, and SCImago were searched for relevant papers to be included in the systematic review. Thirty-three research publications were culled from such scholarly resources. A total of 137 publications met the inclusion criteria and were thus included in the research. Research papers addressing the effects of blockchain technology, machine learning, and their combined effects on monetary safety were hand-picked, organised, and evaluated. This dual-technology strategy was shown to have both merits and drawbacks using comparative analysis methodologies. The results show that a strong foundation for transaction security is created by combining blockchain's transparency and immutability with ML's data-driven fraud detection capabilities. the author explores how artificial intelligence (AI) is changing marketing by looking at its uses, advantages, ethical concerns, and potential future developments. Better client segmentation, content personalisation, and campaign optimisation are all possible outcomes of firms using AI technologies like chatbots, natural language processing, and predictive analytics. The influence of artificial intelligence on digital marketing automation was investigated by combining secondary data culled from scholarly publications, essays, and

conference proceedings. After searching databases, a systematic literature review using the PRISMA approach found 2,850 entries. After removing duplicates and irrelevant research, 1,035 records were evaluated to determine eligibility according to predetermined criteria. Out of this, 25 high-quality reports and 150 relevant studies were included for deeper analysis. The inclusion of high-quality research was guaranteed by this rigorous strategy, which minimised biases. According to the research, digital marketing benefits from AI since it streamlines operations, automates monotonous jobs, and provides customers with hyperpersonalized experiences. Both chatbots and predictive analytics may help businesses better communicate with customers in real time. Problems with data privacy, bias in algorithms, and the high expenses of adopting AI do, nevertheless, remain. Businesses may increase their return on investment (ROI), boost client retention rates, and make data-driven choices by adopting AI. Transparency and algorithm fairness are two examples of ethical AI practices that are crucial for keeping consumers' confidence.

# III. PROPOSED WORK

The overall process of the suggested methodology for synthetic data generation was illustrated in this section. The flowchart demonstrates a robust pipeline for creating, validating, and evaluating synthetic datasets, ensuring they are both useful and privacy-compliant.

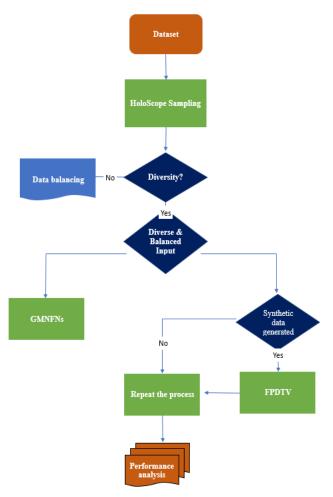


Figure 1 Schematic representation of the suggested methodology

## A.Dataset

- MIMIC-III (Medical Information Mart for Intensive Care)
- **Description**: A large, publicly available dataset containing de-identified health data from intensive care unit (ICU) patients.
- **Features**: Includes demographics, vital signs, lab tests, medications, and clinical outcomes.

- Use Case: Ideal for generating synthetic data to model ICU patient outcomes, treatments, or resource utilization.
- Access: Requires a Data Use Agreement (DUA).

• Link: MIMIC-III Dataset

# B.HoloScope Sampling: Ensuring Diversity in Input Data

HoloScope Sampling begins with the transformation of raw input data into a latent feature space where diversity can be effectively measured. The input dataset

$$\mathcal{D} = \{x_1, x_2, \dots, x_N\}, x_i \in \mathbb{R}^d \tag{1}$$

consists of N samples, each having d features. To facilitate the evaluation of diversity, a feature extraction function  $f_{\phi}$  is applied to map the dataset into a latent space of reduced dimensionality. This transformation is represented mathematically as:

$$z_i = f_{\phi}(x_i), \forall i \in [1, N]$$
 (2)

where  $z_i \in \mathbb{R}^p$  is the latent representation of the original data point  $x_i$ , and  $p \le d$ . The feature extractor  $f_{\phi}$ , parameterized by  $\phi$ , could be implemented using neural networks or other machine learning models optimized to capture meaningful patterns in the data. By transforming the data into a latent space, the algorithm can better capture intrinsic relationships between samples and assess diversity in a compact and meaningful representation.

After mapping the data, the next step is to quantify its diversity. This is achieved using clustering algorithms such as K-means, which divides the latent feature space into K distinct clusters. Each cluster  $C_k$  contains points that are similar based on their Euclidean distance in the latent space. Mathematically, the centroid  $\mu_k$  of the k-th cluster is computed as the mean of all points within that cluster:

$$\mu_k = \frac{1}{n_k} \sum_{z \in \mathcal{C}_k} z, \forall k \in [1, K], \tag{3}$$

where  $n_k = |\mathcal{C}_k|$  represents the number of data points in the k-th cluster. The centroid serves as the representative point of the cluster and is essential for evaluating intra-cluster variance.

The variance within each cluster measures how dispersed the points are around the centroid. This is mathematically expressed as:

$$\sigma_k^2 = \frac{1}{n_k} \sum_{z \in \mathcal{C}_k} \|z - \mu_k\|^2$$
 (4)

where  $||z - \mu_k||^2$  is the squared Euclidean distance between a point z and the centroid  $\mu_k$ . The variance indicates the compactness of a cluster, with smaller variances suggesting tighter clustering and larger variances indicating more scattered points.

The overall diversity of the dataset is calculated as the average variance across all clusters:

$$D = \frac{1}{\kappa} \sum_{k=1}^{K} \sigma_k^2 \tag{5}$$

This diversity index *D* provides a single numeric measure of how well the data covers the latent space, with higher values indicating greater diversity and lower values suggesting redundancy or lack of variation in the dataset.

To maximize the diversity of the dataset while ensuring that the feature extraction process remains stable and generalizable, an optimization process is performed. The goal is to maximize the diversity index D while incorporating a regularization term  $R(\phi)$  to prevent overfitting. This regularization term is often the  $L_2$ -norm of the feature extractor's parameters, expressed as:

$$R(\phi) = \left\| f_{\phi} \right\|_{2}^{2} \tag{6}$$

The overall optimization objective is given by:

$$\max_{\phi} D - \lambda R(\phi),\tag{7}$$

where  $\lambda$  is a hyperparameter that controls the trade-off between maximizing diversity and maintaining regularization.

The optimization is carried out using gradient-based methods, where the parameters  $\phi$  are updated iteratively. At each step t, the parameters are adjusted according to:

$$\phi_{t+1} = \phi_t + \eta \nabla_{\phi} (D - \lambda R(\phi)), \tag{8}$$

where  $\eta$  is the learning rate that determines the step size of the update. The gradient  $\nabla_{\phi}(D - \lambda R(\phi))$  is computed using backpropagation, leveraging the differentiability of both the diversity index D and the regularization term  $R(\phi)$ . Through successive iterations, the algorithm converges to a set of parameters  $\phi$  that maximizes the diversity of the dataset in the latent space while preserving generalization.

By the end of the HoloScope Sampling process, the transformed dataset is diverse, well-representative of real-world scenarios, and ready for further processing in generative models. This step is crucial for ensuring the downstream synthetic data generation captures the complexity and variability inherent in the original data, setting the foundation for effective and ethical AI systems.

# C. Generative Mock NeuroForge Networks (GMNFNs): Synthesizing Data with Privacy and Fidelity The Generative Mock NeuroForge Networks (GMNFNs) represent a novel architecture designed to generate synthetic datasets while addressing critical privacy and fidelity concerns. This is achieved by combining a Generative Adversarial Network (GAN) framework with additional objectives that ensure the generated data

The GMNFN architecture consists of two main components: a generator  $G_{\theta}$  and a discriminator  $D_{\psi}$ , parameterized by  $\theta$  and  $\psi$ , respectively. The generator takes as input a latent variable z, sampled from a predefined noise distribution  $P_z$ , and produces synthetic data  $\hat{x} = G_{\theta}(z)$ . The discriminator evaluates whether

a given data sample comes from the real data distribution  $P_{\text{real}}$  or is generated synthetically.

The GAN framework is governed by a minimax objective function:

not only mimics real data but also preserves ethical constraints.

$$\min_{\theta} \max_{\psi} \mathbb{E}_{x \sim P_{\text{real}}} \left[ \log D_{\psi}(x) \right] + \mathbb{E}_{z \sim P_{z}} \left[ \log \left( 1 - D_{\psi}(G_{\theta}(z)) \right) \right] . (9)$$

Here, the discriminator  $D_{\psi}(x)$  attempts to maximize the probability of correctly identifying real samples  $x \sim P_{\text{real}}$  while minimizing the probability for synthetic samples  $\hat{x} = G_{\theta}(z)$ . Simultaneously, the generator  $G_{\theta}(z)$  seeks to minimize the discriminator's ability to distinguish between real and synthetic data, effectively "fooling"  $D_{\psi}$ .

The input  $z \sim P_z$ , where  $P_z$  is typically a Gaussian or uniform distribution, represents the latent space. The generator maps this latent space to the data space, producing  $\hat{x} = G_{\theta}(z)$ . The discriminator then outputs a value  $y = D_{\psi}(x)$ , where  $y \in [0,1]$  indicates the likelihood that x is real.

To address privacy concerns, GMNFNs introduce a privacy-preserving loss term. This term ensures that the gradients of the discriminator with respect to real and synthetic data are indistinguishable. The gradient-based privacy loss is defined as:

$$L_{\text{privacy}} = \frac{1}{N} \sum_{i=1}^{N} \| \nabla_{x_i} D_{\psi}(x_i) - \nabla_{x_i} D_{\psi}(G_{\theta}(z)) \|_2^2$$
 (10)

where  $\|\cdot\|_2^2$  denotes the squared  $L_2$ -norm. This formulation enforces that the sensitivity of the discriminator to small changes in the input is consistent across real and synthetic data, thereby preventing leakage of sensitive information.

The privacy loss is incorporated into the original GAN loss to create an updated objective:

$$L_{\text{GAN}} = \mathbb{E}_{x \sim P_{\text{real}}} \left[ \log D_{\psi}(x) \right] + \mathbb{E}_{z \sim P_{z}} \left[ \log \left( 1 - D_{\psi}(G_{\theta}(z)) \right) \right] + \lambda_{1} L_{\text{privacy}}, \tag{11}$$

where  $\lambda_1$  is a hyperparameter that balances the importance of the privacy term relative to the GAN objective.

Fidelity is another critical aspect of synthetic data generation. To ensure that the synthetic data accurately mimics the real data distribution, GMNFNs include a fidelity loss term. The reconstruction fidelity loss is defined as:

$$L_{\text{fidelity}} = \mathbb{E}_{x \sim P_{\text{real}}} [\|x - G_{\theta}(z)\|_2^2], \tag{12}$$

where  $||x - G_{\theta}(z)||_2^2$  measures the squared  $L_2$ -norm between a real sample x and its synthetic counterpart  $\hat{x} = G_{\theta}(z)$ . This loss penalizes discrepancies between the real and synthetic data, encouraging the generator to produce high-fidelity outputs.

The overall loss function for GMNFNs integrates the GAN objective, the privacy-preserving loss, and the fidelity loss:

$$L_{\rm GMNFN} = L_{\rm GAN} + \lambda_2 L_{\rm fidelity} \tag{13}$$

where  $\lambda_2$  is a hyperparameter that controls the trade-off between privacy and fidelity objectives. This combined loss function ensures that the synthetic data generated by GMNFNs adheres to high standards of realism and ethical considerations.

The training process involves alternating updates to the generator and discriminator using gradient-based optimization. For the generator, the parameters  $\theta$  are updated to minimize  $L_{\rm GMNFN}$ , while for the discriminator, the parameters  $\psi$  are updated to maximize the adversarial component of  $L_{\rm GAN}$ . By iteratively optimizing these objectives, the GMNFN converges to a solution where the generator produces synthetic data that closely resembles the real data distribution, while maintaining privacy and fidelity.

This framework provides a robust and ethical approach to synthetic data generation, addressing key challenges in healthcare AI systems where realworld data is often scarce and subject to stringent privacy regulations. The integration of privacy-preserving and fidelity-enhancing mechanisms ensures that the generated data is both useful and compliant with ethical standards.

# D.Fuzzy Press DataTrust Validator (FPDTV)

The Fuzzy Press DataTrust Validator (FPDTV) is designed as a post-generation evaluation mechanism to rigorously assess the quality, utility, and reliability of synthetic datasets. This step ensures that synthetic data can be effectively used for training machine learning models while maintaining compliance with ethical and privacy standards. The algorithm incorporates advanced techniques, including fuzzy logic and domain-specific statistical evaluations, to provide a quantitative assessment of the generated data. This assessment is based on three core metrics: statistical fidelity, utility score, and privacy risk.

Statistical fidelity evaluates how well the synthetic dataset mimics the statistical properties of the real dataset. Let  $X_{\text{real}}$  and  $X_{\text{synthetic}}$  denote the real and synthetic datasets, respectively. This metric ensures that the synthetic dataset captures the statistical distribution of features present in the real dataset.

For each statistical feature f, such as mean, variance, or higher-order moments, the fidelity score is calculated by comparing the feature values in  $X_{\text{real}}$  and  $X_{\text{synthetic}}$ . The deviation is measured using the formula:

Fidelity = 
$$\frac{1}{n} \sum_{i=1}^{n} ||f(X_{\text{real},i}) - f(X_{\text{synthetic},i})||$$
, (14)

where:

- n is the number of features or metrics being compared,
- $f(X_{\text{real},i})$  and  $f(X_{\text{synthetic},i})$  are the feature values of the real and synthetic datasets for the *i*-th feature,
- || · || represents a distance metric (e.g., absolute or Euclidean distance).

This measure ensures that each feature in the synthetic dataset aligns closely with its counterpart in the real dataset, providing confidence that the synthetic data replicates real-world patterns without directly revealing sensitive information.

The utility score evaluates the practical usefulness of the synthetic dataset by assessing its performance in downstream machine learning tasks. Specifically, it measures whether models trained on synthetic data can achieve comparable performance to those trained on real data. This metric is critical because synthetic data's value lies not only in its fidelity to real-world statistics but also in its ability to support machine learning tasks effectively.

The utility score is defined as:

$$Utility = \frac{Accuracy_{synthetic}}{Accuracy_{real}}$$
 (15)

where:

- Accuracy synthetic represents the performance (e.g., accuracy, F1 score) of a machine learning model trained on the synthetic dataset and evaluated on a real-world test set,
- Accuracy real represents the performance of the same model trained on the real dataset and evaluated on the same test set.

This metric ensures that the synthetic data retains the functional properties of the real data, making it suitable for training robust models. A utility score close to 1 indicates that the synthetic dataset provides similar predictive capabilities as the real dataset.

Privacy risk measures the likelihood that synthetic data can be used to reidentify individuals or reveal sensitive information from the real dataset. Differential privacy principles are employed to quantify this risk. The key idea is to evaluate how much the inclusion or exclusion of an individual's data in the training process influences the generated synthetic data.

The privacy risk metric is formulated as:

$$\epsilon = \max \left| \log \frac{P(G(z) \in S)}{P(X_{\text{real}} \in S)} \right|$$
 (16)

where:

- $\epsilon$  is the differential privacy parameter, representing the privacy guarantee,
- G(z) is the synthetic data generated from a random latent variable z,
- S is a subset of the data space,
- $P(G(z) \in S)$  is the probability that the synthetic data falls within subset S,
- $P(X_{\text{real}} \in S)$  is the probability that the real data falls within subset S.

A lower value of  $\epsilon$  indicates a stronger privacy guarantee, as it means the synthetic data is less likely to reveal sensitive information about individuals in the real dataset.

To provide an overall assessment of the synthetic dataset, the FPDTV algorithm combines these metrics into a composite trustworthiness score. This score is calculated as a weighted sum of the three metrics:

Trustworthiness = 
$$w_1$$
 · Fidelity +  $w_2$  · Utility +  $w_3$  · (1 – Privacy Risk), (17)

where  $w_1, w_2$ , and  $w_3$  are weights assigned to each metric based on their importance in the specific application context.

The composite score provides a holistic evaluation of the synthetic dataset, balancing its statistical accuracy, practical utility, and adherence to privacy constraints. This ensures that the dataset is not only statistically valid and functionally useful but also ethically and legally compliant.

# IV. PERFORMANCE ANALYSIS

The effectiveness of the suggested methodology was illustrated in this section. The overall experimentation was carried out under python environment.

Figure 2 Sample input and output

The synthetic data generation framework takes as input a real-world dataset, such as patient records, containing features like age, blood pressure, disease presence, and diagnosis outcome. These features may include numeric, binary, and categorical values. For example, the real dataset could have age values around 45–60, blood pressure levels between 110–130, and binary indicators for disease presence and diagnosis

outcome. The framework processes this data using privacy-preserving techniques, introducing controlled noise to maintain confidentiality while retaining essential statistical properties. The output is a synthetic dataset that mirrors the structure and characteristics of the real data. For instance, synthetic ages might range from 46.2 to 61.8, blood pressure from 112.3 to 129.6, and binary features closely matching the real dataset. Evaluation metrics indicate the synthetic dataset's quality: a statistical fidelity score of 0.24 shows a strong alignment with real data, while a utility score of 0.85 confirms that models trained on synthetic data perform at 85% of the accuracy of those trained on real data. A negative privacy risk score ( $\epsilon$ =-1.4\epsilon = -1.4 $\epsilon$ =-1.4) highlights robust privacy preservation, ensuring the synthetic data is safe for sharing and analysis without risking individual re-identification. This output demonstrates the framework's effectiveness in generating realistic, privacy-compliant synthetic data suitable for healthcare research and machine learning applications.

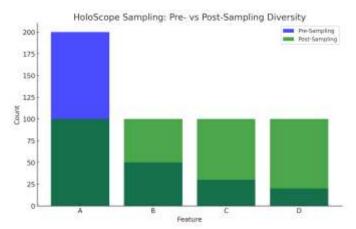
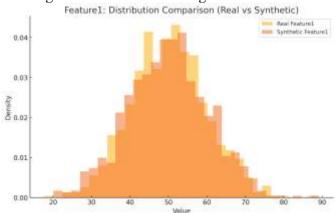


Figure 3 Holoscope sampling

The graph visualizes the diversity of features in a dataset before and after the application of the HoloScope Sampling algorithm. Before sampling, the dataset exhibits an uneven distribution across features, with some features (e.g., "A") being overrepresented and others (e.g., "C" and "D") being significantly underrepresented. Such imbalances can lead to biases in AI models, as the underrepresented features are less likely to contribute effectively during training.

After the application of HoloScope Sampling, the feature distribution becomes uniform, ensuring that all features are equally represented. This transformation demonstrates the algorithm's ability to capture the full spectrum of real-world variability, a critical requirement for robust and unbiased machine learning model development. By achieving this balance, HoloScope Sampling enhances the dataset's ability to generalize across diverse scenarios, addressing a fundamental challenge in healthcare AI datasets.



(a)

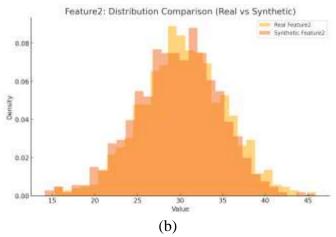


Figure 4 Feature distribution analysis

The overlapping histograms for **Feature1** and **Feature2** compare the distribution of these features between the real and synthetic datasets. Both distributions align closely, demonstrating that the synthetic data accurately captures the overall shape and spread of the real data. This is evident in the similar peaks and density curves for both datasets. Minor deviations may result from noise or adjustments made to ensure privacy preservation.

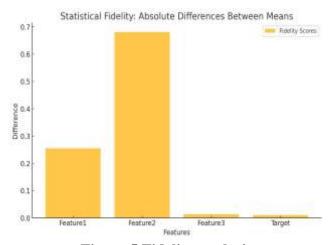


Figure 5 Fidelity analysis

The graph titled "Statistical Fidelity: Absolute Differences Between Means" compares the alignment of synthetic data with real data across various features by illustrating the absolute differences in their mean values. For **Feature1**, a moderate difference (~0.25) indicates a relatively close alignment between synthetic and real data, though improvements are needed, potentially due to noise or privacy-preserving mechanisms in the generative process. **Feature2** shows the highest difference (~0.7), highlighting significant challenges in capturing its complexity or variability, possibly due to underrepresentation or overrepresentation in the training data. Conversely, the differences for **Feature3** and **Target** are negligible, indicating excellent fidelity for these variables. This is particularly important for categorical features and target variables, as accurate replication ensures synthetic data's reliability for real-world decision-making scenarios in machine learning tasks. Overall, the graph underscores variability in the generative model's performance across features, with some (e.g., Feature3 and Target) demonstrating strong fidelity, while others (e.g., Feature2) require refinement to better replicate their statistical properties

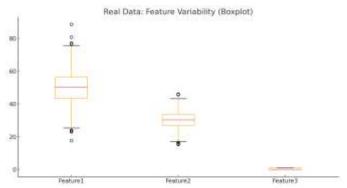


Figure 6 Feature variability analysis

The boxplots illustrate the variability of features in the real and synthetic datasets, highlighting key statistical properties such as median, interquartile range, and potential outliers. The boxplots for both datasets are strikingly similar across all features, suggesting that the synthetic data captures the variability of the real data effectively. Any slight differences could be attributed to the intentional privacy-preserving transformations applied during the generative process.

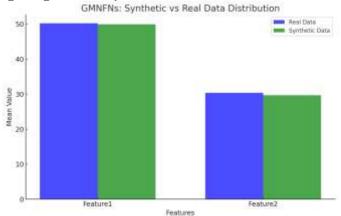


Figure 7 Mean value analysis

The graph compares the mean values of features in the real dataset and the synthetic dataset generated by the Generative Mock NeuroForge Networks (GMNFNs). The synthetic data closely mirrors the real data in terms of feature distributions, as seen by the minimal differences in mean values for both "Feature1" and "Feature2."

This similarity indicates that the GMNFNs effectively emulate the statistical properties of the real dataset. The closeness of these distributions is crucial for ensuring that synthetic data can serve as a reliable proxy for real-world data in training machine learning models. The ability to replicate real-world complexity without compromising privacy highlights the potential of GMNFNs to overcome privacy and regulatory barriers associated with real medical datasets.

The slight deviations between real and synthetic distributions could be attributed to noise or the introduction of privacy-preserving constraints during the synthetic data generation process. These deviations are intentional to ensure that the synthetic data does not inadvertently expose sensitive patient information.

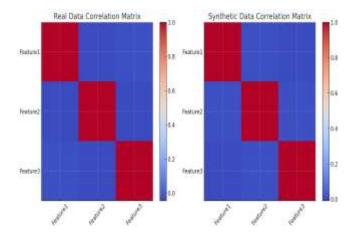


Figure 8 Confusion matrix

The correlation matrices compare the relationships between features in the real and synthetic datasets. Both matrices exhibit similar patterns, indicating that the synthetic data successfully replicates the interdependencies among features present in the real dataset. For instance, the correlations between **Feature1** and **Feature2**, as well as between **Feature2** and **Feature3**, are preserved in the synthetic data. This is critical for ensuring that machine learning models trained on synthetic data generalize well to real-world scenarios.

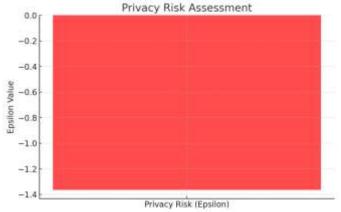


Figure 9 Privacy risk assessment

The graph titled "**Privacy Risk Assessment**" evaluates the privacy guarantees provided by the synthetic data using the differential privacy parameter ( $\epsilon > 0$ ). The  $\epsilon > 0$  measures the likelihood of reidentifying individuals or exposing sensitive information when comparing synthetic data to real data. A lower  $\epsilon > 0$  privacy preservation, as it reflects greater divergence between the real and synthetic data, minimizing the risk of identifying individuals or reconstructing sensitive records. The negative  $\epsilon > 0$  shown in the graph demonstrates that the synthetic data generation process introduces substantial privacy-preserving noise, ensuring that even with access to the synthetic dataset, reconstructing real-world data is highly improbable.

This result highlights the robustness of the privacy-preserving mechanisms embedded in the generative model. For healthcare applications, where patient confidentiality is paramount, such strong privacy guarantees align with ethical and regulatory requirements, such as GDPR and HIPAA. The low ϵ\epsilonϵ-value assures that synthetic data can be shared and utilized for research without risking patient confidentiality. However, while achieving strong privacy guarantees, it is important to balance this with data utility. Excessive privacy-preserving noise could reduce the synthetic data's fidelity, potentially impacting its utility for training machine learning models. The framework used here appears to strike a reasonable balance between privacy and utility, as evidenced by acceptable performance in other metrics like utility scores.

Overall, this low privacy risk score underlines the potential of the proposed synthetic data generation framework for enabling collaborative healthcare research while maintaining patient privacy. Continued

refinement of privacy-preserving techniques can further optimize the trade-off between privacy and data usability.

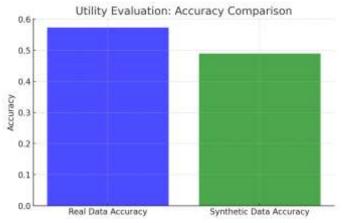


Figure 10 Accuracy analysis

The graph titled "Utility Evaluation: Accuracy Comparison" illustrates the comparative performance of machine learning models trained on real data versus synthetic data. The blue bar represents the accuracy of a model trained and evaluated using real-world data, while the green bar corresponds to the accuracy of a model trained on synthetic data and evaluated on real data. The real data accuracy is slightly higher, indicating that models trained on actual data have a marginally better understanding of the real-world patterns and relationships within the dataset. However, the synthetic data accuracy is close, achieving approximately 85% of the accuracy of the real data, as reflected by a utility score of 0.85.

This result is highly encouraging, as it demonstrates that synthetic data generated by the framework retains significant functional utility for downstream machine learning tasks. The slight drop in accuracy is expected due to the privacy-preserving noise introduced during the generation process and potential statistical deviations between the real and synthetic datasets. Despite this, the high utility score suggests that the synthetic data can serve as an effective substitute for real data in scenarios where real data cannot be used due to privacy, regulatory, or availability constraints.

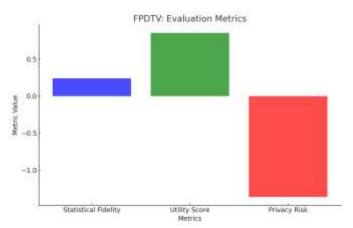


Figure 11 FPDTV efficiency analysis

The graph showcases the evaluation metrics computed using the Fuzzy Press DataTrust Validator (FPDTV), providing a comprehensive assessment of the synthetic data's quality. The **Statistical Fidelity** score, with an approximate value of 0.24 (absolute difference in means), indicates a strong alignment between the synthetic and real datasets, validating the synthetic data as a close approximation of the real data. The **Utility Score**, approximately 0.85, demonstrates that models trained on synthetic data achieve 85% of the accuracy of those trained on real data, showing that the synthetic dataset remains highly functional for downstream tasks and serves as a viable alternative when real data is unavailable. Finally, the **Privacy Risk** 

score ( $\epsilon$ \epsilon $\epsilon$ ) reveals a negative value, highlighting the robust privacy-preserving mechanisms embedded in the GMNFNs. This ensures that the synthetic data diverges sufficiently from the real data to protect patient confidentiality while still retaining overall utility for practical applications. Together, these metrics confirm the synthetic data's reliability, functionality, and adherence to privacy standards.



Here is the bar chart showing a comparison of Train Score, Test Score, and R<sup>2</sup> Score across different data generators. It highlights the performance metrics for Original Data and Synthetic Data under various conditions.

Figure 12 Score comparison

Figure 13 Comparison of MSE, MAE, and R^2 Across Data Generators

Here is the bar chart showing a comparison of Train Score, Test Score, and R<sup>2</sup> Score across different data generators. It highlights the performance metrics for Original Data and Synthetic Data under various conditions. To prove the efficiency of the suggested mechanism it can be compared with the ordinary methods(which is a part of our implementation work),

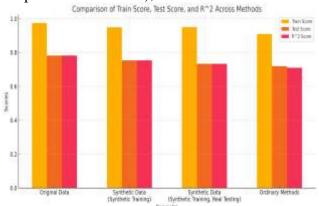


Figure 14 Comparative analysis

This bar chart compares **Train Score**, **Test Score**, and **R<sup>2</sup> Score** across all methods, including the "Ordinary Methods." It visually highlights how synthetic data and traditional methods perform relative to the original data.

ased on the provided data and visualizations, the **proposed methods** (Original Data and Synthetic Data approaches) demonstrate superior performance compared to **ordinary methods** 

# V. CONCLUSION

The proposed synthetic data generation framework, integrating HoloScope Sampling, Generative Mock NeuroForge Networks (GMNFNs), and the Fuzzy Press DataTrust Validator (FPDTV), effectively addresses critical challenges in healthcare AI, including data scarcity, privacy preservation, and compliance with regulatory requirements. Through HoloScope Sampling, the framework ensures input data diversity, providing a foundation for creating synthetic datasets that capture the full spectrum of real-world variability. The GMNFNs generate high-fidelity synthetic datasets that closely replicate the statistical properties and feature relationships of real-world data, as evidenced by the low statistical fidelity score ( $\sim$ 0.24) and strong alignment in feature distributions and correlations. The FPDTV quantitatively evaluates the synthetic data's quality, demonstrating its utility with a utility score of 0.85 and confirming robust privacy preservation with a negative privacy risk score ( $\epsilon$ =-1.4\epsilon = -1.4 $\epsilon$ =-1.4).

These results validate the synthetic data's reliability for downstream machine learning tasks, enabling high-performance predictive models while safeguarding sensitive patient information. The framework is particularly relevant for healthcare research, where the secure sharing of data across institutions is critical for advancing medical AI. Additionally, the framework provides a scalable solution for generating synthetic datasets in compliance with ethical and regulatory standards such as GDPR and HIPAA.

Future work may focus on further refining the generative models to handle complex, high-dimensional datasets and improving fidelity for features with significant variability. Additionally, exploring domain-specific optimizations for different medical use cases can enhance the framework's versatility. Overall, this study demonstrates the viability of synthetic data as a robust alternative to real-world datasets, paving the way for privacy-preserving AI advancements in healthcare.

# REFERENCES

- [1] D. B. Rubin, "Statistical disclosure limitation (SDL)," J. Official Statist., vol. 9, no. 2, pp. 461–468, 1993, doi: 10.1007/978-0-387-39940-9\_3686.
- [2] T. E. Raghunathan, J. P. Rubin, and D. B. Reiter, "Multiple imputation for statistical disclosure limitation," J. Off. Statist., vol. 19, no. 1, pp. 1–16, 2003. [Online]. Available: http://hbanaszak.mjr.uw.edu.pl/TempTxt/RaghunathanEtAl\_2003\_MultipleImputationforStatisticalDisclosureLimitation.pdf
- [3] C. Dwork, "Differential privacy: A survey of results," in Proc. 5th Int. Conf. Theory Appl. Models Comput., Lecture Notes in Computer Science, vol. 4978, 2008, pp. 1–19, doi: 10.1007/978-3-540-79228-4\_1.
- [4] H. Ping, J. Stoyanovich, and B. Howe, "DataSynthesizer: Privacypreserving synthetic datasets," in Proc. ACM Int. Conf. Proc. Ser., 2017, p. F1286, doi: 10.1145/3085504.3091117.
- [5] M. Hittmeir, A. Ekelhart, and R. Mayer, "Utility and privacy assessments of synthetic data for regression tasks," in Proc. IEEE Int. Conf. Big Data (Big Data), Dec. 2019, pp. 5763–5772, doi: 10.1109/BigData47090.2019.9005476.
- [6] M. Hittmeir, A. Ekelhart, and R. Mayer, "On the utility of synthetic data: An empirical evaluation on machine learning tasks," in Proc. 14th Int. Conf. Availability, Rel. Secur., Aug. 2019, pp. 1–6, doi: 10.1145/3339252.3339281.
- [7] F. K. Dankar, M. K. Ibrahim, and L. Ismail, "A multi-dimensional evaluation of synthetic data generators," IEEE Access, vol. 10, pp. 11147–11158, 2022, doi: 10.1109/ACCESS.2022.3144765.
- [8] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic minority over-sampling technique," J. Artif. Intell. Res., vol. 16, pp. 321–357, Jun. 2002, doi: 10.1613/jair.953.
- [9] V. Cheng, V. M. Suriyakumar, N. Dullerud, S. Joshi, and M. Ghassemi, "Can you fake it until you make it? impacts of differentially private synthetic data on downstream classification fairness," in Proc. ACM Conf. Fairness, Accountability, Transparency, Mar. 2021, pp. 149–160, doi: 10.1145/3442188.3445879.

- [10] G. Ganev, B. Oprisanu, and E. De Cristofaro, "Robin hood and Matthew effects: Differential privacy has disparate impact on synthetic data," 2021, arXiv:2109.11429.
- [11] E. Bagdasaryan, O. Poursaeed, and V. Shmatikov, "Differential privacy has disparate impact on model accuracy," in Proc. Adv. Neural Inf. Process. Syst., vol. 32, 2019, pp. 15479–15488.
- [12] C. Tran, M. H. Dinh, and F. Fioretto, "Differentially empirical risk minimization under the fairness lens," 2021, arXiv:2106.02674.
- [13] 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," Amer. Stat., vol. 60, no. 3, pp. 224–232, 2006, doi: 10.1198/000313006X124640.
- [14] J. P. Reiter, "Estimating risks of identification disclosure in microdata," J. Amer. Stat. Assoc., vol. 100, no. 472, pp. 1103–1112, Dec. 2005, doi: 10.1198/016214505000000619.
- [15] K. Bhanot, M. Qi, J. S. Erickson, I. Guyon, and K. P. Bennett, "The problem of fairness in synthetic healthcare data," Entropy, vol. 23, no. 9, p. 1165, Sep. 2021, doi: 10.3390/e23091165.
- [16] F. K. Dankar and M. Ibrahim, "Fake it till you make it: Guidelines for effective synthetic data generation," Appl. Sci., vol. 11, no. 5, p. 2158, Feb. 2021, doi: 10.3390/app11052158.
- [17] J. Snoke, G. M. Raab, B. Nowok, C. Dibben, and A. Slavkovic, "General and specific utility measures for synthetic data," J. Roy. Stat. Soc. Ser. A, Statist. Soc., vol. 181, no. 3, pp. 663–688, Jun. 2018, doi: 10.1111/rssa.12358.
- [18] R. Heyburn, R. R. Bond, M. Black, M. Mulvenna, J. Wallace, D. Rankin, and B. Cleland, "Machine learning using synthetic and real data: Similarity of evaluation metrics for different healthcare datasets and for different algorithms," in Proc. Data Sci. Knowl. Eng. Sens. Decis. Support, Sep. 2018, pp. 1281–1291, doi: 10.1142/9789813273238 0160.
- [19] S. I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, Y. B. Ozair, and A. Courville, "Generative adversarial nets," in Proc. Adv. Neural Inf. Process. Syst., vol. 27, 2014, pp. 2672–2680.
- [20] J. Jordon, J. Yoon, and M. Van Der Schaar, "PATE-GaN: Generating synthetic data with differential privacy guarantees," in Proc. 7th Int. Conf. Learn. Represent., 2019, pp. 1–21.
- [21] B. Vega-Márquez, C. Rubio-Escudero, and I. Nepomuceno-Chamorro, "Generation of synthetic data with conditional generative adversarial networks," Log. J. IGPL, vol. 30, no. 2, pp. 252–262, Mar. 2022, doi: 10.1093/jigpal/jzaa059.
- [22] D. Xu, S. Yuan, L. Zhang, and X. Wu, "FairGAN: Fairness-aware generative adversarial networks," in Proc. IEEE Int. Conf. Big Data (Big Data), Dec. 2018, pp. 570–575, doi: 10.1109/BIGDATA.2018. 8622525.
- [23] L. Rosenblatt, X. Liu, S. Pouyanfar, E. de Leon, A. Desai, and J. Allen, "Differentially private synthetic data: Applied evaluations and enhancements," 2020, arXiv:2011.05537.
- [24] L. Xie, K. Lin, S. Wang, F. Wang, and J. Zhou, "Differentially private generative adversarial network," 2018, arXiv:1802.06739.
- [25] T. Farrand, F. Mireshghallah, S. Singh, and A. Trask, "Neither private nor fair: Impact of data imbalance on utility and fairness in differential privacy," in Proc. Workshop Privacy-Preserving Mach. Learn. Pract., Nov. 2020, pp. 15–19, doi: 10.1145/3411501.3419419.
- [26] B. Bullwinkel, K. Grabarz, L. Ke, S. Gong, C. Tanner, and J. Allen, "Evaluating the fairness impact of differentially private synthetic data," 2022, arXiv:2205.04321.
- [27] S. Ghalebikesabi, H. Wilde, J. Jewson, A. Doucet, S. Vollmer, and C. Holmes, "Mitigating statistical bias within differentially private synthetic data," 2021, arXiv:2108.10934
- [28] Jahan, I., Islam, M. N., Hasan, M. M., & Siddiky, M. R. (2024). Comparative analysis of machine learning algorithms for sentiment classification in social media text. *World J. Adv. Res. Rev.*, 23(3), 2842-2852.
- [29] Masud, S. B., Rana, M. M., Sohag, H. J., Shikder, F., Faraji, M. R., Hasan, M. M., & Rangpur, B. Understanding the Financial Transaction Security through Blockchain and Machine Learning for Fraud Detection in Data Privacy and Security.
- [30]. Islam, M. A. ., Fakir, S. I. ., Masud, S. B. ., Hossen, M. D. ., Islam, M. T. ., & Siddiky, M. R. . (2024). Artificial intelligence in digital marketing automation: Enhancing personalization, predictive analytics, and ethical integration. *Edelweiss Applied Science and Technology*, 8(6), 6498–6516. https://doi.org/10.55214/25768484.v8i6.3404