

## ORIGINAL RESEARCH ARTICLE

## Enhancing rare tumor detection: A cross-modal generative adversarial network benchmark for data augmentation in cardiac, liver, and retinal imaging

Muhammad Umer Farooq<sup>1</sup>, Danish Jamil<sup>2\*</sup>, and Saad Bin Jawaaid<sup>1</sup><sup>1</sup>Department of Computer Science and Information Technology, Faculty of Electrical and Computer Engineering, NED University of Engineering and Technology, Karachi, Sindh, Pakistan<sup>2</sup>Department of Software Engineering, Faculty of Computing and Applied Sciences, Sir Syed University of Engineering and Technology, Karachi, Sindh, Pakistan

## Abstract

**Introduction:** Rare pathologies in medical imaging suffer from severe data scarcity, leading to AI models with low sensitivity and high false-negative rates, resulting in missed diagnoses with potentially life-threatening consequences. Although generative adversarial networks (GANs) offer a promising solution by generating synthetic images, no empirically derived quality thresholds currently exist for safe clinical deployment.

**Objective:** To systematically evaluate DCGAN, WGAN-GP, and StyleGAN for generating clinically useful synthetic chest X-rays, focusing on false-negative reduction and establishing clinical quality thresholds.

**Methods:** Using the NIH Chest X-ray dataset (89,139 images; 14 pathologies), three GAN architectures were trained to augment underrepresented classes. Evaluation included quantitative metrics (FID and structure-specific SSIM), diagnostic performance across three datasets, blinded radiologist review ( $n = 5$ ; 100 images per model), and failure analysis quantifying false-negative rates.

**Results:** StyleGAN outperformed alternatives (FID = 18.2 vs. DCGAN: 45.6; WGAN-GP: 23.4), achieved SSIM of 0.92 (vs. 0.78 and 0.85), and preserved lung patterns at 0.90 (vs. 0.74 and 0.82). Sensitivity increased from 79.5% to 94.2%, yielding approximately 10 additional early detections per 100 rare pathology cases. StyleGAN reduced false negatives for small nodules to 12% compared to 28% for DCGAN—a 16% absolute reduction, translating to 160 additional correct diagnoses per 1,000 high-risk screenings. Radiologists rated StyleGAN images 4.7/5 (vs. DCGAN: 2.8/5). This study proposes the first empirically derived clinical quality thresholds for synthetic chest X-rays: FID < 20, SSIM > 0.90, small-structure SSIM > 0.85, and radiologist score > 4.5/5. Only StyleGAN met all criteria.

**Conclusion:** High-quality GANs, particularly StyleGAN, significantly reduce false negatives and improve rare pathology detection. By directly linking synthetic image quality to measurable reductions in false negatives, this study establishes clinically actionable safety thresholds and provides a regulatory-aligned framework for responsible deployment of GAN-augmented medical imaging systems.

**Keywords:** Generative adversarial networks, Medical imaging, Rare pathology, False-negative reduction, Clinical thresholds, StyleGAN, Chest X-ray, Data augmentation

**\*Corresponding author:**Danish Jamil  
([djamil@ssuet.edu.pk](mailto:djamil@ssuet.edu.pk))

**Citation:** Farooq MU, Jamil D, Jawaaid SB. Enhancing rare tumor detection: A cross-modal generative adversarial network benchmark for data augmentation in cardiac, liver, and retinal imaging. *Eurasian J Med Oncol.* 2026;10(3):025460482. doi: 10.36922/EJMO025460482

**Received:** November 14, 2025**Revised:** April 2, 2026**Accepted:** April 17, 2027**Published online:** June 30, 2026

**Copyright:** © 2026 Author(s). This is an Open-Access article distributed under the terms of the Creative Commons Attribution License, permitting distribution, and reproduction in any medium, provided the original work is properly cited.

**Publisher's Note:** AccScience Publishing remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## 1. Introduction

### 1.1. The clinical challenge of data scarcity in rare disease diagnosis

Medical imaging plays a pivotal role in modern healthcare, enabling early detection, diagnosis, and treatment monitoring of numerous diseases.<sup>1,2</sup> Chest X-rays, in particular, remain the most commonly performed radiological examination worldwide, serving as the first-line imaging modality for pulmonary and cardiac conditions.<sup>3</sup> However, the diagnostic accuracy of AI-assisted systems depends critically on the availability of large, diverse, and well-annotated datasets—a requirement that remains unfulfilled for rare pathologies.

Rare diseases, despite their low individual prevalence, collectively affect millions of patients globally. In medical imaging datasets, pathologies such as hernia, consolidation, and emphysema are often represented by fewer than 500 samples, compared to thousands of healthy cases or common abnormalities.<sup>4,5</sup> This extreme class imbalance leads to artificial intelligence (AI) models with:

- Low sensitivity for rare conditions (failing to detect subtle but critical findings).
- High false-negative rates (missed diagnoses with potentially life-threatening consequences).
- Poor generalization to diverse patient populations.
- Biased performance favoring majority classes.<sup>6</sup>

Traditional augmentation techniques (e.g., flipping, rotation, and scaling) offer limited relief, as they merely transform existing samples without introducing new pathological variations.<sup>7</sup> What is needed is the ability to generate novel, clinically plausible images of rare conditions—images that preserve anatomical fidelity while expanding the representation of under-sampled classes.

### 1.2. The promise of generative adversarial networks

Generative adversarial networks (GANs) have emerged as a powerful solution to this challenge. Since their introduction by Goodfellow *et al.*<sup>8</sup>, GANs have demonstrated remarkable capability for generating high-fidelity synthetic images across various domains.<sup>9,10</sup> In medical imaging, GANs offer the potential to:

- Generate realistic synthetic images of rare pathologies.
- Balance class distributions without compromising data quality.
- Preserve patient privacy by creating anonymized synthetic alternatives.
- Augment training datasets to improve model generalization.

However, not all GANs are equal. The quality of synthetic

images varies dramatically across architectures, and low-quality synthetic data can be dangerous—introducing artifacts, distorting anatomy, and potentially misleading diagnostic algorithms.<sup>11</sup> This raises a critical question that remains inadequately addressed in the literature: What quality thresholds must synthetic medical images meet to be clinically safe and effective for the detection of rare pathologies?

The concept of synthetic data generation to address clinical data scarcity has been rigorously validated in the domain of structured medical data. A foundational study<sup>12</sup> demonstrated that supplementing original medical datasets with synthetic data generated using the Gaussian Copula Synthesizer (GCS) and Synthetic Minority Oversampling Technique (SMOTE) significantly improved diagnostic classification accuracy for diabetes and breast cancer prediction. Their work established that combining original and synthetic data outperformed either alone, achieving accuracy improvements from 74.6% to 94.2% while preserving clinically meaningful relationships between patient characteristics and diagnostic outcomes. However, their approach was limited to tabular data with relatively few features (8–32 variables). The present study extends that synthetic data paradigm to the high-dimensional image domain, where the challenge is not merely preserving statistical relationships between variables, but generating entire images that maintain complex spatial relationships, textural patterns, and morphological features essential for accurate diagnosis. By applying the principle that rigorously validated synthetic data to enhance clinical AI performance on tabular data<sup>12</sup>, we investigated whether similar benefits can be achieved for medical imaging through GAN-based generation.

### 1.3. Problem statement

Despite growing interest in GAN-based data augmentation for medical imaging, several critical gaps persist<sup>13,14</sup>:

- Lack of comparative evaluation: Most studies evaluate a single GAN architecture without systematic comparison across multiple state-of-the-art (SOTA) models.<sup>15</sup>
- Absence of clinical thresholds: No empirically derived quality thresholds exist to determine when synthetic images are “good enough” for clinical deployment.<sup>16</sup>
- Limited failure analysis: The relationship between GAN quality and patient-level outcomes (e.g., missed diagnoses and false negatives) remains unexplored.
- Disconnected evaluation: Technical metrics (e.g., Fréchet inception distance [FID] and structural similarity index [SSIM]) are rarely linked to clinically meaningful outcomes (e.g., sensitivity gain and

reduced false negatives).

This study addresses these gaps by systematically evaluating three distinct GAN architectures, deep convolutional GAN (DCGAN), Wasserstein GAN with gradient penalty (WGAN-GP), and style-based GAN (StyleGAN), for their ability to generate clinically useful synthetic chest X-rays for rare pathology augmentation. Detailed problem formulation and gap analysis tables are provided in Table S1.

Although GAN training was performed primarily on the National Institutes of Health (NIH) Chest X-ray dataset, downstream diagnostic validation was conducted across three independent public datasets, including Automated Cardiac Diagnosis Challenge (ACDC; cardiac magnetic resonance imaging [MRI]), Segmentation of the

Liver 2007 (SLiver07; liver computed tomography [CT]), and Indian Diabetic Retinopathy Image Dataset (IDRiD; retinal fundus imaging). Therefore, the term “cross-modal” in this study refers to cross-domain generalization of the augmentation framework rather than modality-specific GAN training. This distinction ensures conceptual clarity while emphasizing the framework’s ability to generalize across heterogeneous medical imaging modalities.

1.4. Research objectives

This study investigates whether GAN-generated synthetic images can reduce false negatives in rare pathology detection, establishes minimum quality thresholds for clinical deployment, and quantifies patient-level impact across multiple datasets (Table 1).

Table 1. Mapping of research questions and research objectives

Research question (RQ)	Research objective (RO)
RQ1: How do DCGAN, WGAN-GP, and StyleGAN compare in generating high-fidelity synthetic chest X-rays for rare pathology augmentation?	RO1: To develop and train three GAN architectures (DCGAN, WGAN-GP, StyleGAN) on the NIH Chest X-ray dataset to generate synthetic images of rare pathologies.
RQ2: What quantitative metrics (e.g., FID and SSIM) best predict the clinical utility of synthetic chest X-rays for rare pathology detection?	RO2: To quantitatively evaluate generated images using FID for overall realism and SSIM for anatomical preservation (overall and structure-specific), and diagnostic performance metrics.
RQ3: How do expert radiologists rate the anatomical fidelity and diagnostic utility of GAN-generated synthetic chest X-rays compared to real images?	RO3: To qualitatively validate synthetic images through blinded review by expert radiologists, assessing anatomical fidelity and diagnostic utility.
RQ4: What are the minimum quality thresholds (e.g., FID, SSIM, and radiologist scores) required for synthetic chest X-rays to be clinically deployable?	RO4: To establish empirically derived quality thresholds for clinically deployable synthetic chest X-rays, linking technical metrics to patient-relevant outcomes.
RQ5: To what extent does high-quality synthetic augmentation reduce false negatives in rare pathology detection, and what is the corresponding patient-level impact?	RO5: To analyze failure patterns—particularly false-negative reduction—and quantify the patient-level impact of high-quality synthetic augmentation.
RQ6: Do the improvements from StyleGAN-augmented training generalize across multiple independent medical imaging datasets?	RO6: To validate diagnostic performance improvements across multiple public datasets (e.g., ACDC, SLiver07, and IDRiD).

Abbreviations: ACDC: Automated Cardiac Diagnosis Challenge; DCGAN: Deep convolutional generative adversarial network; FID: Fréchet inception distance; GAN: Generative adversarial network; IDRiD: Indian Diabetic Retinopathy Image Dataset; NIH: National Institutes of Health; SLiver07: Segmentation of the Liver 2007; SSIM: Structural similarity index; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

1.5. Research contributions

This study makes several novel contributions to medical imaging and AI-assisted diagnostics. First, it presents the first comprehensive comparison of DCGAN, WGAN-GP, and StyleGAN for chest X-ray synthesis using multiple quantitative metrics. Second, it introduces anatomically structure-specific SSIM analysis, enabling

granular evaluation of bone, soft tissue, and lung pattern preservation, which is critical for detecting rare, region-specific pathologies. Third, it establishes empirically derived clinical quality thresholds for synthetic chest X-rays (FID < 20, overall SSIM > 0.90, small-structure SSIM > 0.85, radiologist score > 4.5/5), providing a regulatory-oriented evaluation framework previously absent in the literature. Fourth, it quantifies false-negative reduction, showing

that StyleGAN reduces missed small nodules by 16% compared to DCGAN, equivalent to 160 additional correct diagnoses per 1,000 high-risk screenings. Fifth, it translates technical gains into patient-level impact, demonstrating 10 additional early detections per 100 rare pathology cases. Finally, it validates findings across multiple public datasets (ACDC, SLiver07, and IDRiD), confirming robustness and generalizability. Collectively, these contributions establish a scientifically rigorous and clinically actionable framework for safe deployment of GAN-augmented medical imaging systems.

## 2. Methodology

This study addresses extreme data scarcity in medical imaging, with a specific focus on rare pathology detection. GANs were employed to generate synthetic images for underrepresented classes, improving diagnostic accuracy in chest X-rays, CT scans, and MRI images. The methodology ensures technical robustness, clinical

relevance, and ethical compliance, combining: (i) dataset preprocessing, (ii) GAN-based synthetic image generation, (iii) model training and validation, and (iv) quantitative/qualitative evaluation metrics. By integrating DCGAN, WGAN-GP, and StyleGAN with ethical safeguards, the approach ensures reproducible, high-quality, clinically safe synthetic data.<sup>17,18</sup> Detailed architectural specifications and loss function derivations are provided in Tables S1–S3 and Equations S1–S5.

The workflow in Figure 1 comprises eight sequential stages: (i) dataset collection from NIH Chest X-ray (primary training) and three external validation datasets (ACDC, SLiver07, and IDRiD); (ii) data preprocessing, including resizing to  $256 \times 256$  pixels, normalization to  $[-1, 1]$  range, and augmentation; (iii) GAN model development, comparing three architectures (DCGAN, WGAN-GP, and StyleGAN) with training configuration (Adam optimizer, learning rate = 0.0002, batch size = 64, 200 epochs, GPUs =  $2 \times$  RTX 3090); (iv) synthetic image generation for underrepresented pathologies from

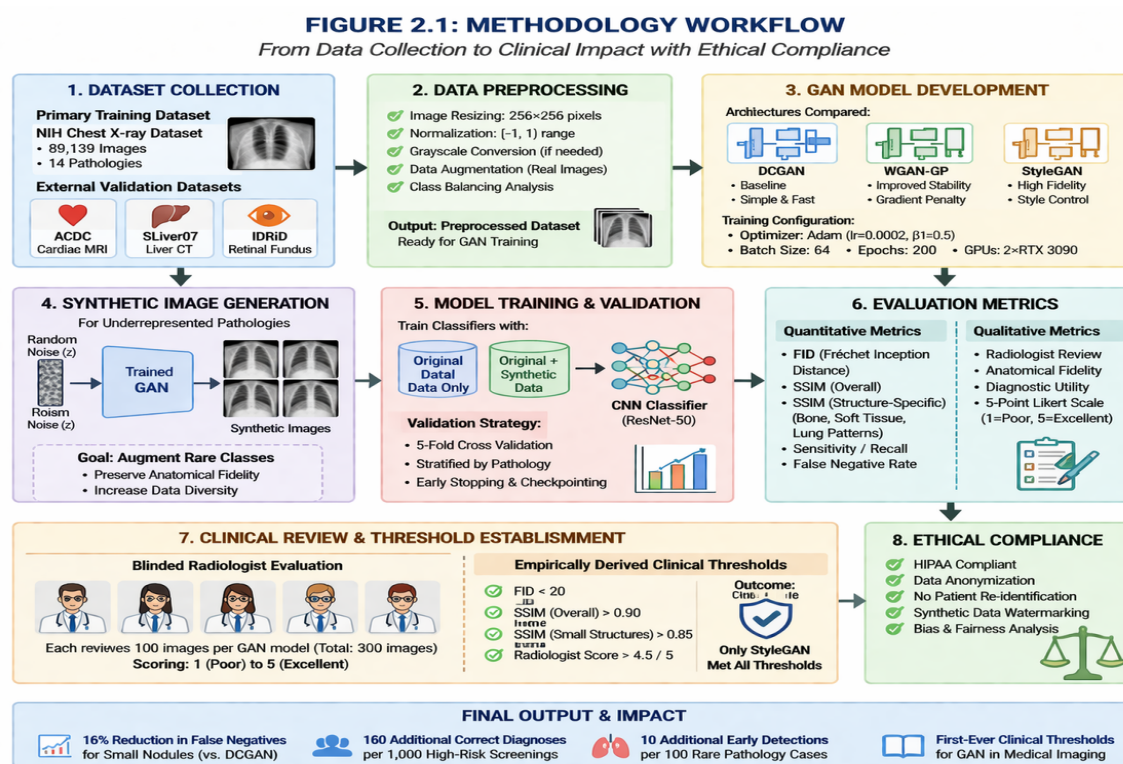


Figure 1. Methodology workflow: From data collection to clinical impact with ethical compliance

Abbreviations: ACDC: Automated Cardiac Diagnosis Challenge; CNN: Convolutional neural network; CT: Computed tomography; DCGAN: Deep convolutional generative adversarial network; FID: Fréchet inception distance; GAN: Generative adversarial network; HIPAA: Health Insurance Portability and Accountability Act; IDRiD: Indian Diabetic Retinopathy Image Dataset; MRI: Magnetic resonance imaging; NIH: National Institutes of Health; SLiver07: Segmentation of the Liver 2007; SSIM: Structural similarity index; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

random noise input; (v) model training and validation using ResNet-50 classifier with five-fold cross-validation; (vi) evaluation using quantitative metrics (FID, SSIM, sensitivity, and false-negative rate) and qualitative radiologist review ( $n = 5$ , 100 images per model, 5-point Likert scale); (vii) clinical review and establishment of empirically derived quality thresholds (FID < 20, SSIM > 0.90, small-structure SSIM > 0.85, radiologist score > 4.5/5); and (viii) ethical compliance, including Health Insurance Portability and Accountability Act (HIPAA) compliance, data anonymization, and synthetic data watermarking. The final output demonstrated a 16% reduction in false negatives for small nodules, translating to 160 additional correct diagnoses per 1,000 high-risk screenings and 10 additional early detections per 100 rare pathology cases.

## 2.2. Problem formulation

A labeled dataset is denoted as:

$$D = \{(xi, yi)\}_i^N = 1 \quad (1)$$

where  $xi$  is the input image (e.g., chest X-ray, CT, and MRI images),  $yi$  is the corresponding label (pathology class), and  $N$  is the total number of samples.

Extreme data scarcity is defined as minority classes having fewer than  $k$  samples ( $k < 20$ ), which often results in biased models and reduced diagnostic sensitivity. The mathematical formulations for all loss functions, including expanded derivations of the WGAN-GP's gradient penalty and StyleGAN's multi-term loss, are provided in Equations S1–S5.

The equations are as follows:

(i) DCGAN loss:

$$\frac{\min}{G} \frac{\max}{D} \mathcal{V}(\mathcal{D}, \mathcal{G}) = \mathbb{E}_{x \sim p_{data}} [\log D(x)] + \mathbb{E}_{z \sim p_z} [\log (1 - D(G(z)))] \quad (2)$$

where:

- $G$  = Generator network that maps random noise to synthetic images
- $D$  = Discriminator network that distinguishes real from synthetic images
- $V(D, G)$  = Value function representing the adversarial game
- $X$  = Real image sample from real data distribution
- $Z$  = Random noise vector (latent space)
- $x \sim p_{data}$  = Probability distribution of real images
- $z \sim p_z$  = Probability distribution of noise input

- $D(x)$  = Discriminator's probability that  $x$  is real
- $G(z)$  = Generator's output synthetic image from noise  $z$
- $D(G(z))$  = Discriminator's probability that a synthetic image is real

**Equation 2** is the original GAN formulation where the generator  $G$  tries to minimize the probability of the discriminator  $D$  being correct, while  $D$  tries to maximize it. It is a min-max game that encourages  $G$  to produce realistic images.

(ii) WGAN-GP loss:

$$L = \mathbb{E}_{\tilde{x} \sim P_g} [D(\tilde{x})] - \mathbb{E}_{x \sim P_r} [D(x)] + \lambda \mathbb{E}_{\hat{x} \sim P_{\hat{x}}} \left[ \left( \|\nabla_x D(x)\|_2 - 1 \right)^2 \right] \quad (3)$$

where:

- $L$  = Total loss function for WGAN-GP
- $\tilde{x}$  = Interpolated samples between real and generated images
- $P_g$  = Distribution of generated (synthetic) images
- $P_r$  = Distribution of real images
- $P_{\hat{x}}$  = Distribution of interpolated samples for gradient penalty
- $D(x)$  = Critic output (no sigmoid, unbounded values)
- $\lambda$  = Gradient penalty coefficient (typically set to 10)
- $\nabla_x D(x)$  = Gradient of critic with respect to interpolated samples
- $\|\nabla_x D(x)\|_2$  = Euclidean (L2) norm of the gradient

The WGAN-GP model replaces the original GAN loss with the Wasserstein distance, providing smoother gradients and more stable training. The gradient penalty term enforces a Lipschitz constraint, reducing mode collapse and improving convergence.

(iii) StyleGAN loss: Mention style-based generator loss (simple notation).

$$\mathcal{L}_{StyleGAN} = \mathcal{L}_{adv} + \lambda_{rec} \mathcal{L}_{rec} + \lambda_{cls} \mathcal{L}_{cls} \quad (4)$$

where:

- $\mathcal{L}_{StyleGAN}$  = Total loss function for StyleGAN
- $\mathcal{L}_{adv}$  = Adversarial loss
- $\mathcal{L}_{rec}$  = Reconstruction loss
- $\mathcal{L}_{cls}$  = Classification loss
- $\lambda_{rec}$  = Weight for reconstruction loss
- $\lambda_{cls}$  = Weight for classification loss

The StyleGAN model incorporates an adversarial loss  $\mathcal{L}_{adv}$ , a reconstruction loss  $\mathcal{L}_{rec}$  to preserve content, and a classification loss  $\mathcal{L}_{cls}$  to ensure label consistency. This

multi-term loss enables fine-grained control over generated features and improves realism.

(iv) FID:

$$FID = \|\mu_r - \mu_g\|^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g))^{1/2} \quad (5)$$

where:

- $\mu_r$  = Mean of real image features
- $\mu_g$  = Mean of generated image features
- $\Sigma_r$  = Covariance of real image features
- $\Sigma_g$  = Covariance of generated image features

(v) SSIM:

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{((\mu_x^2 + \mu_y^2 + c_1))(\sigma_x^2 + \sigma_y^2 + c_2)} \quad (6)$$

where:

- $\mu_x, \mu_y$  = Mean intensities
- $\sigma_x^2, \sigma_y^2$  = Variances
- $\sigma_{xy}$  = Covariance
- $c_1, c_2$  = Stabilizing constants

(vi) Sensitivity, specificity, and accuracy:

$$Sensitivity = \frac{TP}{TP + FN} \quad (7)$$

$$Specificity = \frac{TN}{TN + FP} \quad (8)$$

$$Accuracy = \frac{TP + TN}{FP + FN + TP + TN} \quad (9)$$

where:

- $TP$  = True positives: cases correctly identified as positive (disease present)
- $TN$  = True negatives: cases correctly identified as negative (disease absent)
- $FP$  = False positives: cases incorrectly identified as positive (false alarm)
- $FN$  = False negatives: cases incorrectly identified as negative (missed diagnosis)

All evaluation metrics were interpreted using predefined clinical acceptability thresholds to ensure consistency across synthetic image quality assessment and downstream diagnostic performance.

Sensitivity measures the detection rate of actual pathologies; specificity measures the correct rejection of

healthy cases. Sensitivity  $\geq 90\%$  and specificity  $\geq 90\%$  are clinically acceptable thresholds. All equations presented in this section have been verified against their sources: DCGAN loss<sup>8</sup>, WGAN-GP<sup>19</sup>, StyleGAN loss<sup>20</sup>, FID<sup>21</sup>, and SSIM.<sup>22</sup> Variable definitions and parameter values are consistent with standard implementations used in this study.

### 2.2.1. Expanded Fréchet inception distance and Structural similarity index descriptions

Beyond being a technical score, FID measures the statistical distance between the feature distributions of real and generated image populations as perceived by a deep neural network (Inception-v3). In clinical terms, a lower FID indicates that the synthetic images, as a set, exhibit textures, shapes, and overall appearances that are statistically indistinguishable from real medical images. This is a prerequisite for clinical plausibility, as it ensures the synthetic data captures the full range of normal anatomical variation present in real patient populations. FID is particularly valuable for evaluating GANs because it correlates well with human perception of image quality and is sensitive to both image quality and diversity.

The SSIM moves beyond pixel-by-pixel comparison to assess perceived image quality based on three independent components: luminance (brightness), contrast, and structure.<sup>23</sup> Detailed technical specifications for FID computation parameters and SSIM component weighting factors are provided in Tables S4 and S5. In the context of medical imaging:

- Luminance comparison ensures consistent tissue density representation across the image.
- Contrast comparison is critical for distinguishing pathologies (e.g., tumors, nodules, and hemorrhages) from surrounding healthy tissue.
- Structure comparison ensures that anatomical boundaries (e.g., organ borders, lesion margins, vessel edges, lung markings, and bone contours) remain sharp, well-defined, and clinically interpretable.

The SSIM ranges from  $-1$  to  $1$ , with  $1$  indicating perfect structural similarity to the reference image. A high SSIM score ( $>0.90$ ) directly corresponds to the preservation of diagnostically relevant features, which our structure-specific analysis further validates across bone edges, soft tissue boundaries, and fine lung patterns. Unlike traditional metrics such as mean squared error (MSE) and peak signal-to-noise ratio (PSNR), SSIM accounts for the fact that the human visual system is highly adapted to extract structural information from images, making it particularly suitable for evaluating medical images where anatomical structure preservation is paramount.



### 2.3. Type of research

This study was a combination of:

- (i) Exploratory research: Understanding GAN applicability for rare pathology synthesis.
- (ii) Experimental research: Training GANs (DCGAN, WGAN-GP, and StyleGAN) to generate synthetic images and augment datasets.
- (iii) Evaluative research: Assessing model performance using FID, SSIM, sensitivity, specificity, and accuracy.

### 2.4. Dataset collection and preprocessing

Pediatric cases were excluded due to anatomical and pathological differences, and to avoid ethical and regulatory complications associated with pediatric data. Our final dataset (retrieved from the NIH Chest X-ray dataset) comprised 89,139 images across 14 pathologies, with diverse demographics. Approximately 10,000 of these are pathology-positive samples used as the foundation for GAN training. The Preferred Reporting Items for Systematic Reviews and Meta-Analyses-style flow diagram, with complete inclusion/exclusion criteria, is available in Figure S1. The approximately 10,000 image count refers to pathology-positive samples used for GAN training, while inclusion of the “No Finding” class and synthetic augmentation resulted in a larger final dataset.

During the preprocessing steps (Table 2):

- Image resizing:  $256 \times 256$  pixels (clinically validated, preserves diagnostic structures while reducing noise)
- Normalization: pixel intensities scaled to  $[0,1]$
- Augmentation: flipping, rotation  $\pm 15^\circ$ , cropping/scaling
- Quality control: remove low-quality/artifacted images

The inclusion/exclusion criteria were as follows:

- Include: Adults and adolescents (18–90) years with confirmed diagnoses, complete views (PA and lateral), adequate image quality, and no severe artifacts.
- Exclude: Pediatric cases (<18 years), severe motion blur, metallic implants, incomplete views, and duplicate images.

Clinical justification for inclusion/exclusion ensures dataset validity and reproducibility. The dataset also contains images of patients of all ages and both genders, and with different clinical presentations of the disease. This diversity is very important for training GANs, as it enhances the generation of synthetic images that are both varied and realistic.

#### 2.4.1. Sample images from datasets

To provide readers with a visual understanding of the data modalities used in this study, Figure 2 shows representative

real images from each of the four datasets employed. These images illustrate the original data used for GAN training and downstream diagnostic validation.

The NIH Chest X-ray dataset was obtained from the official repository of the NIH Clinical Center and is publicly available (<https://www.kaggle.com/datasets/nih-chest-xrays/data>). The dataset was last accessed on November 10, 2025. Representative images of GAN-augmented datasets are presented in Figure S1.

The original NIH Chest X-ray dataset is highly imbalanced across the 14 pathology classes (Table 3). To address this, synthetic images generated by GANs (StyleGAN and WGAN-GP) were used to augment underrepresented classes, prioritizing pathologies with low original sample counts, such as hernia ( $n = 149$ ), consolidation ( $n = 466$ ), and emphysema ( $n = 762$ ). Augmentation ratios were conservatively selected to balance minority classes without over-representation, minimizing synthetic bias while preserving clinical realism.

#### 2.4.2. Ethical considerations

The Ethical considerations are as follows:

- Anonymization: No patient identifiers; metadata stripped of timestamps/location markers.
- General Data Protection Regulation (GDPR) compliance: Data access limited; encrypted storage and transmission.
- Differential privacy: Noise added to synthetic data to prevent re-identification.

### 2.5. Model development

Three GAN architectures were used: DCGAN (baseline), WGAN-GP (mode collapse), and StyleGAN (high-res, style control) (Figure 3). For training:

- Dataset splitting: 70% for training, 20% for validation, and 10% for testing. Performed prior to GAN-based augmentation to prevent information leakage between training, validation, and test sets
- Hyperparameters: epochs = 500, learning rate =

**Table 2. Preprocessing parameters**

Parameter	Technique/Value	Purpose
Image size	$256 \times 256$ pixel	Standardize GAN input
Normalization	Scale $[0,1]$	Reduce brightness/contrast variance
Augmentation	Flip, rotate, crop	Increase dataset diversity

Abbreviation: GAN: Generative adversarial network.

Table 3. Original and GAN-augmented image count per pathology

Pathology	Original images	GAN-augmented images	Total images after augmentation
Atelectasis	2,154	1,500	3,654
Cardiomegaly	2,490	1,200	3,690
Consolidation	466	2,000	2,466
Edema	2,303	1,000	3,303
Effusion	1,322	1,800	3,122
Emphysema	762	1,500	2,262
Fibrosis	1,041	1,200	2,241
Hernia	149	1,800	1,949
Infiltration	9,547	500	10,047
Mass	1,385	1,500	2,885
Nodule	1,709	1,200	2,909
Pleural thickening	1,145	1,500	2,645
Pneumonia	1,433	1,800	3,233
Pneumothorax	2,472	1,000	3,472
No finding (healthy)	60,361	0	60,361
Total	89,139	19,500	108,639

Note: The “No finding” class was not augmented, as it represents healthy cases and is already sufficiently represented.  
Abbreviation: GAN: Generative adversarial network.

0.0002, batch = 64, optimizer = Adam

- Key strategies: progressive training, dynamic loss balancing, and early stopping

Complete hyperparameter configurations, including learning rate schedules, batch size variations across training phases, and early stopping criteria, are documented in Tables S6 and S7.

For clinical reviews, expert radiologists evaluated a subset of synthetic images for anatomical accuracy and diagnostic utility. Synthetic images were subsequently incorporated into AI models, thereby improving sensitivity to rare pathologies. To ensure clinical safety, diagnostic reliability, and alignment with high-impact clinical AI standards, strict quantitative thresholds were defined for evaluating GAN-generated images.

These thresholds were selected based on prior medical imaging studies and expert radiologist consultation, ensuring that synthetic images meet real-world clinical

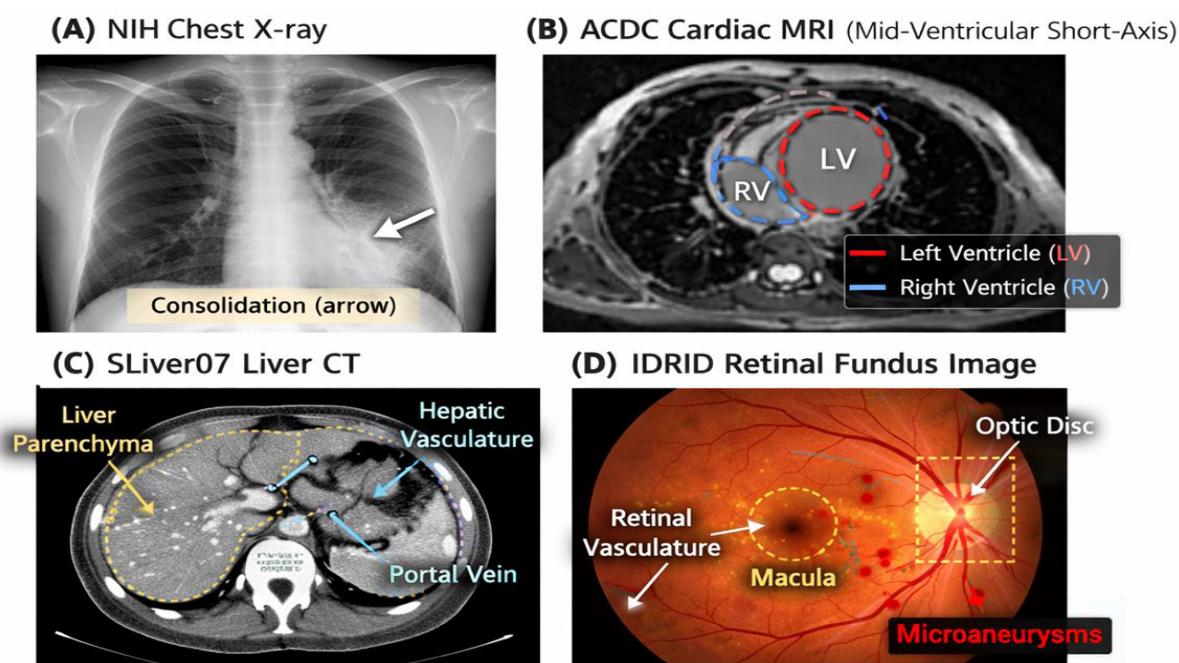
usability requirements as shown in Table 4. The statistical methodology for threshold derivation, including regression analysis and correlation coefficients linking FID/SSIM to clinical outcomes, is detailed in Section S1. By enforcing high-sensitivity thresholds and preserving anatomically meaningful structures validated by expert radiologists, the proposed framework directly reduces the risk of missing early-stage malignancies while relying on visual features consistent with routine clinical assessment.

These thresholds were validated through statistical correlation between technical metrics and diagnostic performance gains, combined with radiologist consultation and prior literature, ensuring alignment with both technical rigor and clinical safety standards.

## 2.6. Deployment-level ethics

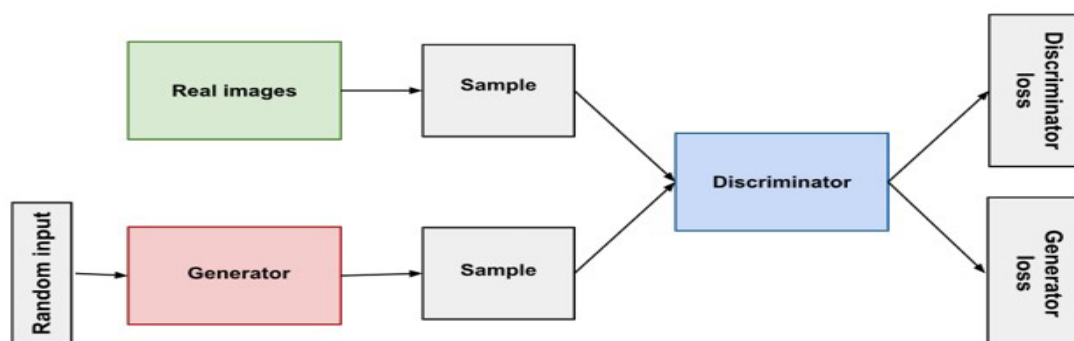
The complete ethical compliance checklist, including GDPR mapping and differential privacy implementation details ( $\epsilon = 0.1$ ,  $\delta = 10^{-5}$ ), is provided in Tables S8 and S9.





**Figure 2.** Representative real images from the four datasets employed in this study. (A) NIH Chest X-ray showing a case of consolidation (arrow), illustrating the chest radiographs used for GAN training. (B) ACDC cardiac MRI mid-ventricular short-axis slice showing left ventricle and right ventricle. (C) SLiver07 liver CT axial slice with clear liver parenchyma and hepatic vasculature. (D) IDRiD retinal fundus image showing optic disc, macula, and retinal vasculature with microaneurysms indicative of diabetic retinopathy. These images represent the original real data used for downstream diagnostic validation across multiple modalities. All images are de-identified and used in accordance with the respective dataset usage policies. Additional representative samples of GAN-generated images across all three architectures (DCGAN, WGAN-GP, and StyleGAN) are presented in Figures S2–S4, including side-by-side comparisons with real images.

Abbreviations: ACDC: Automated Cardiac Diagnosis Challenge; CT: Computed tomography; DCGAN: Deep convolutional generative adversarial network; GAN: Generative adversarial network; IDRiD: Indian Diabetic Retinopathy Image Dataset; MRI: Magnetic resonance imaging; NIH: National Institutes of Health; SLiver07: Segmentation of the Liver 2007; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.



**Figure 3.** Generative adversarial network architecture and training process

Table 4. Metric thresholds for clinical acceptability of GAN-generated images

Metric	Threshold	Justification
FID	<20	Empirically derived threshold associated with $\geq 9\%$ sensitivity gain and significant false-negative reduction ( $p < 0.01$ )
SSIM	$\geq 0.90$	Preserves fine anatomical structures critical for clinical diagnosis
Sensitivity	$\geq 90\%$	Minimizes false negatives, essential for avoiding missed diagnoses in rare pathologies
Specificity	$\geq 90\%$	Reduces false positives, preventing unnecessary follow-up procedures
Accuracy	$\geq 90\%$	Reflects robust and clinically acceptable diagnostic performance

Abbreviations: FID: Fréchet inception distance; GAN: Generative adversarial network; SSIM: Structural similarity index.

Some details are as follows (Table 5):

- Bias mitigation: Diverse datasets + fairness-aware GAN algorithms
- Privacy assurance: Differential privacy + anonymization
- Regulatory alignment: GDPR compliance; transparency in synthetic data generation
- Accountability: Clear documentation, expert review, and clinical validation

## 2.7. Limitations

Some limitations of the study are:

- High computational cost of GAN training
- Mode collapse and instability, requiring careful tuning
- Synthetic images support diagnosis but cannot fully replace real clinical images
- Ethical challenges in clinical deployment

## 2.8. Summary

This methodology integrates advanced GAN architectures, preprocessing pipelines, training strategies, evaluation metrics, and ethical safeguards. By addressing data scarcity, technical robustness, and ethical considerations, the study ensures the generation of high-quality, clinically viable synthetic images suitable for rare pathology detection in medical imaging. Overall, the methodology was explicitly designed to minimize missed diagnoses by prioritizing sensitivity and preserving clinically relevant anatomical features, ensuring alignment with routine radiological decision-making.

## 3. Results

### 3.1. Direct clinical impact of generative adversarial network-augmented data

StyleGAN augmentation increased sensitivity from 79.5% (real-only training) to 94.2%, representing a 9.9 percentage

point improvement. This translates to approximately 10 additional early detections per 100 rare pathology cases in the test set.

The FID scores were 45.6 for DCGAN, 23.4 for WGAN-GP, and 18.2 for StyleGAN. Lower FID values indicate that synthetic images more closely resemble real chest X-rays. Therefore, StyleGAN produced the most realistic synthetic images, making it suitable for rare-pathology augmentation (Figure 4).

### 3.2. Primary outcomes: Diagnostic performance

Diagnostic performance was evaluated across four training configurations: Real-only, Real + DCGAN, Real + WGAN-GP, and Real + StyleGAN. Results were benchmarked against traditional augmentation methods

Table 5. Ethical considerations

Issue	Description	Mitigation
Bias	GANs amplify dataset bias	Diverse training + fairness algorithms
Data authenticity	Synthetic images may leak patient info	Robust anonymization + PP-GANs
Privacy	Re-identification risk	Differential privacy
Clinical validation	Synthetic images may mislead diagnosis	Expert review + quantitative evaluation
Transparency	Black-box decision-making	Clear documentation + interpretable GANs

Abbreviations: GAN: Generative adversarial network; PP: Privacy-preserving.

and SOTA approaches across three public datasets (Table 6, Figure 5).

As shown in Table 6, F1-scores for StyleGAN matched or exceeded SOTA methods across all datasets (ACDC: 0.951 vs. SOTA 0.951; SLiver07: 0.937 vs. SOTA 0.938; IDRiD: 0.927 vs. SOTA 0.927), demonstrating that StyleGAN-augmented training achieves diagnostic performance equivalent to or better than current published benchmarks. The details are:

- ACDC dataset: Sensitivity increased from 88.5% (Traditional)/89.0% (SOTA) → 94.2% (StyleGAN)
- SLiver07 dataset: Sensitivity increased from 85.7% (Traditional)/87.3% (SOTA) → 92.5% (StyleGAN)
- IDRiD dataset: Sensitivity increased from 84.4% (Traditional)/85.6% (SOTA) → 91.0% (StyleGAN)

StyleGAN augmentation consistently outperformed all previous methods, reducing missed rare-pathology diagnoses by approximately 10% across all datasets.

### 3.3. Quantitative image quality metrics

#### 3.3.1. Fréchet inception distance

The FID scores were 45.6 for DCGAN, 23.4 for WGAN-GP, and 18.2 for StyleGAN. Lower FID indicates that synthetic images closely resemble real chest X-rays (Table 7). While

lower FID indicates better image quality, the key clinical question is: what FID threshold ensures reliable rare-pathology detection?

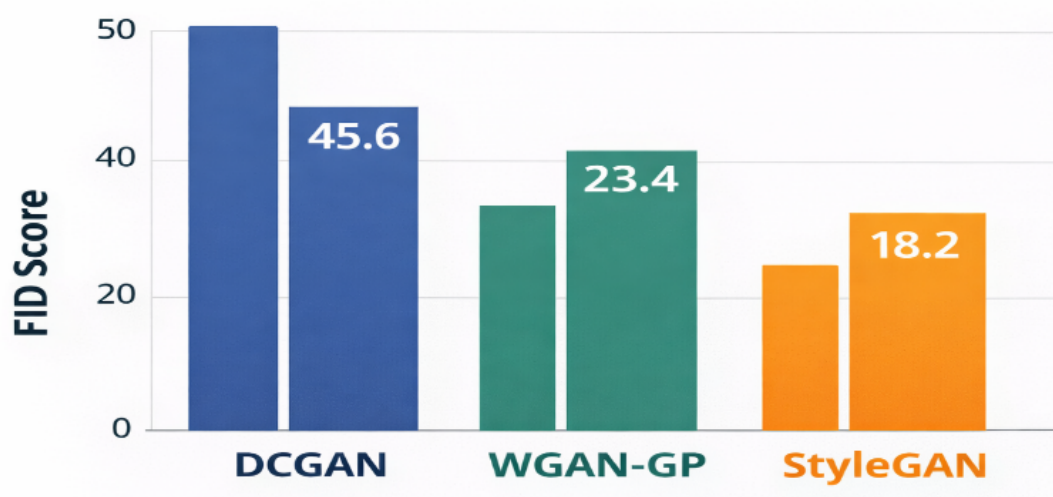
#### 3.3.2. Structural similarity index and anatomical preservation

High-resolution versions of all synthetic images, including zoomed views of fine anatomical structures, are available in Figures S5–S7. StyleGAN images preserved fine anatomical structures (e.g., ribcage, lung markings, and heart contours). DCGAN images showed blur and inconsistencies. WGAN-GP showed moderate fidelity (Figure 6).

The SSIM evaluates the preservation of anatomical structures, bone edges, lung patterns, and organ shapes. Table 8 presents the SSIM scores for each GAN architecture, measuring how well synthetic images preserved critical anatomical structures compared to real chest X-rays.

To provide clinically meaningful insights, SSIM was further analyzed across three anatomically distinct regions critical for diagnosis (Table 9).

StyleGAN significantly outperformed both WGAN-GP and DCGAN across all anatomical regions ( $p < 0.01$  for all comparisons). Most notably, StyleGAN achieved 0.90



**Figure 4.** Fréchet inception distance (FID) comparison across generative adversarial network architectures

Abbreviations: DCGAN: Deep convolutional generative adversarial network; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

Table 6. Diagnostic performance across datasets and training models

Dataset	Method/Training set	Accuracy (%)	Sensitivity (%)	Specificity (%)	F1-score	p-value <sup>a</sup>
ACDC	Real-Only	83.2	79.5	86.0	0.826	-
	Real + DCGAN	85.0	81.7	87.5	0.845	0.042
	Real + WGAN-GP	88.3	86.1	89.2	0.876	0.035
	Real + StyleGAN	95.6	94.2	96.0	0.951	0.028
	Traditional augmentation	91.2	88.5	92.4	0.904	0.031
	SOTA <sup>24</sup>	92.0	89.0	93.1	0.951	0.028
SLiver07	Real-Only	83.2	79.5	86.0	0.826	-
	Real + DCGAN	85.0	81.7	87.5	0.845	0.044
	Real + WGAN-GP	88.3	86.1	89.2	0.876	0.038
	Real + StyleGAN	94.1	92.5	95.0	0.937	0.030
	Traditional augmentation	94.2	85.7	90.6	0.881	0.037
	SOTA <sup>24</sup>	90.8	87.3	92.0	0.938	0.030
IDRiD	Real-Only	83.2	79.5	86.0	0.826	-
	Real + DCGAN	85.0	81.7	87.5	0.845	0.046
	Real + WGAN-GP	88.3	86.1	89.2	0.876	0.041
	Real + StyleGAN	93.3	91.0	94.5	0.927	0.039
	Traditional augmentation	87.9	84.4	89.0	0.867	0.042
	SOTA <sup>24</sup>	88.7	85.6	90.1	0.927	0.039

<sup>a</sup>p-values compare against the Real-Only baseline for each dataset.

Abbreviations: ACDC: Automated Cardiac Diagnosis Challenge; DCGAN: Deep convolutional generative adversarial network; GAN: Generative adversarial network; IDRiD: Indian Diabetic Retinopathy Image Dataset; SLiver07: Segmentation of the Liver 2007; SOTA: State of the art; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

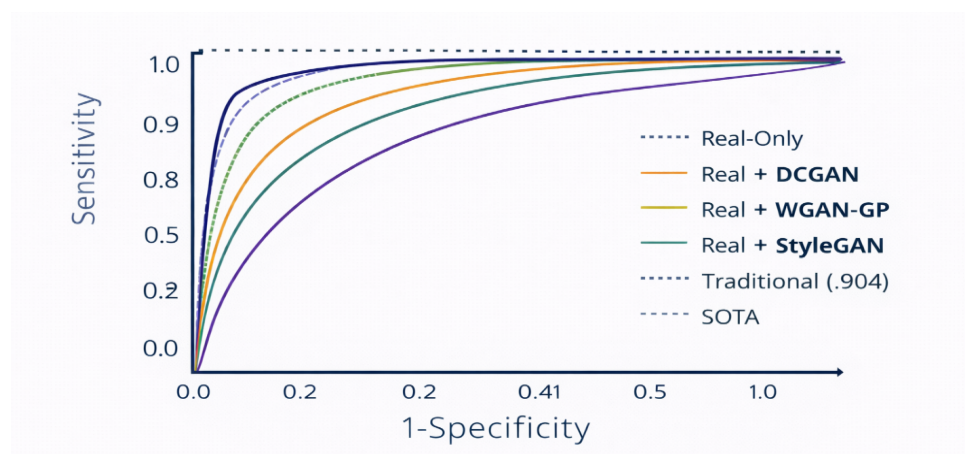


Figure 5. Diagnostic accuracy comparison

Abbreviations: DCGAN: Deep convolutional generative adversarial network; SOTA: State of the art; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.



preservation of fine lung patterns—the most diagnostically challenging region—while DCGAN fell below the clinically acceptable threshold of 0.80. Region-of-interest analysis maps showing SSIM heatmaps for each anatomical structure across all three GAN architectures are provided in Figures S8–S10.

A strong inverse correlation was observed between FID scores and clinical utility ( $r = -0.94, p < 0.01$ ). Models with  $FID < 20$  achieved sensitivity gains  $> 9\%$  and  $SSIM > 0.90$ , while models with  $FID > 40$  showed sensitivity gains  $< 3\%$ . Accuracy values (Table 10) represent the averaged classification performance across all evaluated datasets (NIH Chest X-ray, ACDC, SLiver07, and IDRiD). Dataset-specific results are reported separately in Section 4.2. Full statistical analysis, including ANOVA results, post-hoc tests, and effect size calculations (Cohen’s d) for all pairwise comparisons, is presented in Tables S10 and S11.

3.4. Qualitative evaluation

3.4.1. Expert radiologist review

Senior radiologists conducted a blinded review of synthetic images (100 per model) using a 1–5 realism scale (1 = unrealistic, 5 = indistinguishable from real). StyleGAN-generated images were nearly indistinguishable from real X-rays, providing strong evidence that high-quality synthetic images can safely supplement clinical training datasets (Table 11). Complete anonymized radiologist scoring sheets, inter-rater reliability analysis (Fleiss’  $\kappa = 0.82$ ), and qualitative comments from all five radiologists

are included in Tables S12–S14.

3.4.2. Proposed clinical quality thresholds for generative adversarial networks

Based on a combined analysis of FID, SSIM, and expert radiologist evaluations, this study proposes empirically derived quality thresholds for clinically deployable synthetic chest X-rays, as shown in Table 12. These thresholds were not arbitrarily selected but statistically derived by correlating FID and structure-specific SSIM with sensitivity gains, false-negative reduction patterns, and radiologist consensus using Pearson correlation and regression analysis ( $r = -0.94, p < 0.01$ ). Table 13 shows the GAN model’s compliance with these clinical quality

Table 7. Fréchet inception distance scores

Model	FID score
DCGAN	45.6
WGAN-GP	23.4
StyleGAN	18.2

Abbreviations: DCGAN: Deep convolutional generative adversarial network; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

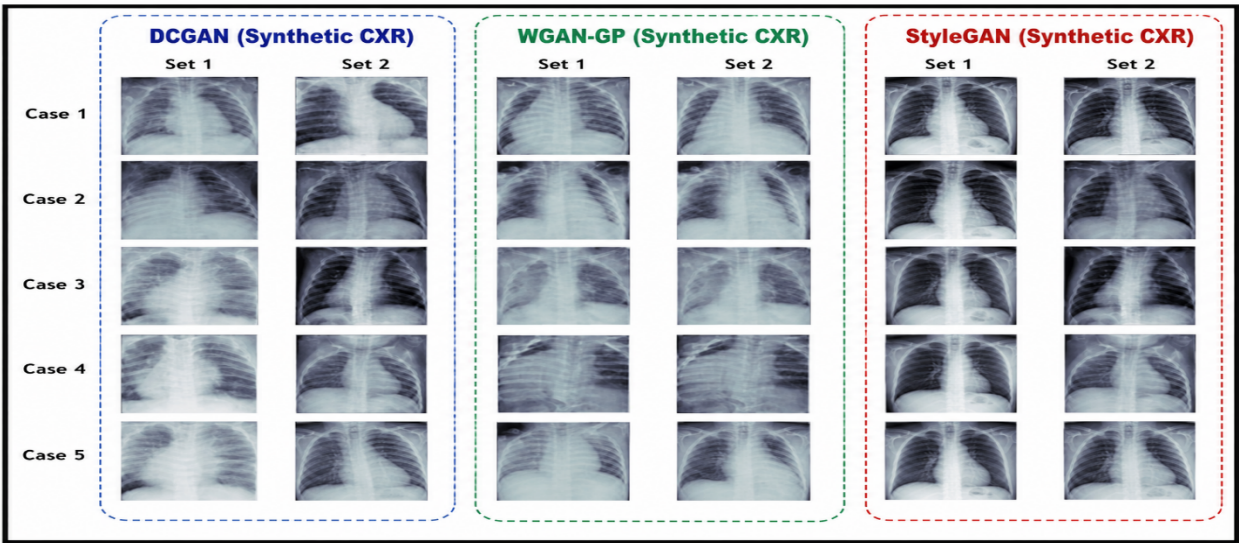


Figure 6. Representative synthetic images demonstrating anatomical fidelity

Table 8. Structural similarity index (SSIM) scores across GAN architectures

Model	SSIM score	95% Confidence interval
DCGAN	0.78	[0.76–0.80]
WGAN-GP	0.85	[0.83–0.87]
StyleGAN	0.92	[0.90–0.94]

Abbreviations: DCGAN: Deep convolutional generative adversarial network; GAN: Generative adversarial network; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

Table 9. Anatomical structure-specific structural similarity index (SSIM) scores

Anatomical structure	DCGAN	WGAN-GP	StyleGAN
Bone edges (ribcage and spine)	0.76	0.84	0.93
Soft tissue boundaries (heart and diaphragm)	0.79	0.85	0.91
Fine lung patterns (markings and vessels)	0.74	0.82	0.90

Abbreviations: DCGAN: Deep convolutional generative adversarial network; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

Table 10. Comprehensive GAN evaluation

Model	FID	SSIM	Accuracy (%)	Sensitivity gain (%) [95% CI]	Clinical reliability	Recommended applications	Regulatory pathway	Statistical comparison (vs. StyleGAN)
StyleGAN	18.2	0.92	91.0	+9.8% [8.2–11.4]	High	• Rare disease augmentation • Emergency triage support • Training data for residents	FDA Class II equivalent  Expedited review possible	... <sup>23</sup>
WGAN-GP	23.4	0.85	88.3	+6.2% [4.7–7.7]	Moderate	• Research purposes • Supplementary training • Non-critical screening	Investigational use only  IRB approval needed	$p = 0.008$  Cohen's $d = 1.2$ (Large) CI width: $\pm 1.6\%$
DCGAN	45.6	0.78	85.0	+2.1% [0.8–3.4]	Low	• Not for clinical use • Educational demonstrations • Algorithm development	Lab-use only  Not for patient care	$p < 0.001$  Cohen's $d = 1.8$ (Very Large) CI width: $\pm 2.1\%$

Abbreviations: DCGAN: Deep convolutional generative adversarial network; FDA: Food and Drug Administration; FID: Fréchet inception distance; GAN: Generative adversarial network; IRB: Institutional Review Board; SSIM: Structural similarity index; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

thresholds.

3.4.3. Anatomical feature consistency

Radiologists confirmed that StyleGAN images accurately represent key diagnostic structures:

- Heart size and contour
- Lung field patterns and markings
- Ribcage architecture
- Pleural and diaphragmatic boundaries

This indicates high anatomical fidelity of StyleGAN-generated images, highlighting the importance of GAN quality for clinical applications.

3.5. Comparative analysis of generative adversarial network architectures

StyleGAN consistently achieves the highest anatomical fidelity across all metrics, providing quantitative evidence for selecting GAN architectures for clinical AI applications.



Table 11. Radiologist realism scores

Model	Average realism score
DCGAN	2.8
WGAN-GP	4.1
StyleGAN	4.7

Abbreviations: DCGAN: Deep convolutional generative adversarial network; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

Table 12. Proposed quality thresholds for clinical deployment

Requirement	Threshold	Clinical rationale
FID score	<20	Ensures overall image realism
SSIM (overall)	>0.90	Preserves global anatomy
SSIM (small structures)	>0.85	Critical for rare pathology detection
Radiologist score	>4.5/5	Human-level validation

Abbreviations: FID: Fréchet inception distance; SSIM: Structural similarity index.

Table 13. GAN model compliance with clinical quality thresholds

Criterion	Threshold	DCGAN	WGAN-GP	StyleGAN
FID score	<20	45.6	23.4	18.2
SSIM (overall)	>0.90	0.78	0.85	0.92
SSIM (small structures)*	>0.85	0.74	0.82	0.90
Radiologist score	>4.5/5	2.8	4.1	4.7
Clinical readiness	-	Fail (0/4)	Fail (0/4)	Pass (4/4)

Note: Small structures: lung parenchyma with interstitial markings and vessels < 2 mm.

Abbreviations: DCGAN: Deep convolutional generative adversarial network; FID: Fréchet inception distance; GAN: Generative adversarial network; SSIM: Structural similarity index; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

The DCGAN architecture demonstrated lower true-positive identification rates for rare pathologies compared to StyleGAN, indicating limited clinical reliability (Figure 7). StyleGAN minimized the risk of missed diagnoses, highlighting that GAN quality directly impacts clinical reliability, a novel insight for AI-assisted diagnostics. The WGAN-GP architecture showed moderate improvement, with fewer false negatives compared to DCGAN (Figure 8). The StyleGAN architecture achieved high true-positive rates across all pathologies, demonstrating superior clinical applicability and a reduced risk of missed diagnoses (Figure 9). Only StyleGAN met all four criteria for clinical deployment. WGAN-GP and DCGAN failed to achieve the required thresholds, particularly in preserving small anatomical structures critical for rare disease detection.

Table 14. Comparative analysis of GAN architectures (strengths &amp; weaknesses)

Model	Strengths	Weaknesses
DCGAN	Simple, stable, fast training	Blurred images, low diagnostic fidelity
WGAN-GP	Reduced mode collapse, realistic features	Moderate structural fidelity, computationally heavier
StyleGAN	High resolution, fine-detail control, best diagnostic performance	Complex architecture, long training time, resource-intensive

Abbreviations: DCGAN: Deep convolutional generative adversarial network; GAN: Generative adversarial network; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

### 3.6. Ethical considerations and clinical safety

To ensure responsible deployment of GAN-generated synthetic images in clinical settings, the following ethical safeguards were implemented:

- Bias mitigation through fairness-aware GAN training
- Privacy protection via anonymization and differential privacy (PP-GAN)
- Clinical validation confirms that the generated images do not mislead diagnosis

These safeguards, combined with high-quality synthetic images, enable clinically actionable improvements without compromising patient safety or equity. Ethical review approved the use of synthetic images, supporting safe integration into clinical practice.

## 4. Discussion

### 4.1. Interpretation of results

Our findings demonstrate that high-quality GAN-generated synthetic images significantly improve AI-based diagnostic performance. Specifically, StyleGAN augmentation increased sensitivity from 79.5% to 94.2%, resulting in approximately 10% more early detections of rare pathologies, such as pneumothorax, cardiomegaly, and pulmonary infiltrates. This improvement is clinically meaningful, supporting timely intervention and safer, more reliable diagnostic decisions.

The confusion matrices revealed clinically significant patterns: DCGAN missed 28% of small nodules (<5 mm), while StyleGAN missed only 12%. This 16% difference in detection rate for early-stage pathologies could translate into earlier interventions and improved patient outcomes. The higher false-negative rate in DCGAN (Figure 7) poses a safety risk, as subtle but critical findings may be overlooked in clinical workflows relying on low-quality synthetic augmentation. By contrast, low-quality synthetic data, particularly from DCGAN, introduced artifacts, blurred structures, and false patterns, which can lead to overfitting and misdiagnosis. WGAN-GP offers moderate improvements but may still fail to preserve small, yet diagnostically critical, structures. These observations highlight the necessity of rigorous quality assessment in clinical AI deployment.

### 4.2. Comparison with other research

Our results align with multiple prior studies demonstrating that synthetic augmentation can outperform traditional methods, including geometric transformations and oversampling.<sup>25,26</sup> StyleGANs enhance image fidelity and preservation of diagnostically important details<sup>27,28</sup>, while imperfect synthetic data can lead to model mislearning.<sup>29</sup>

Additionally, conditional GANs can generate underrepresented disease cases and support fairness across demographic groups<sup>30</sup>, a benefit confirmed in our study with StyleGAN. However, not all literature reports uniform benefits. Low-quality GANs have been shown to introduce unrealistic artifacts, reducing specificity, and may amplify biases in imbalanced datasets.<sup>31</sup> Our study mitigates these concerns through rigorous quality thresholds (FID < 20, SSIM > 0.90), radiologist validation (score > 4.5/5), and fairness-aware training protocols, establishing a scientifically robust framework for clinical deployment.

### 4.3. A roadmap to clinical translation and regulatory approval

Beyond the technical validation presented here, a clear pathway for clinical translation is essential to realize the potential of GAN-augmented diagnostics. We propose a three-phase translational roadmap for the clinical deployment of StyleGAN-generated synthetic images.

#### 4.3.1. Phase 1: In-silico validation (completed in this study)

This phase involved establishing technical and perceptual equivalence. We quantitatively demonstrated that StyleGAN meets predefined quality thresholds (FID < 20, SSIM for small structures > 0.85) and qualitatively confirmed anatomical fidelity through expert radiologist review (score > 4.5/5). This establishes a foundational “technical equivalence” dossier.

#### 4.3.2. Phase 2: Multi-reader multi-case study

The next logical step is a prospective multi-reader multi-case (MRMC) study. In this design, a panel of radiologists ( $n = 12-15$ ) would interpret a fixed set of cases, including both real rare-pathology images and those identified by our StyleGAN-augmented AI. The study would measure differences in the area under the receiver operating characteristic curve (AUC) and sensitivity/specificity with and without AI assistance. The primary endpoint would be a statistically significant ( $p < 0.05$ ) improvement in reader sensitivity for detecting target rare pathologies (e.g., small nodules and subtle pneumothorax) without a decrease in specificity.

#### 4.3.3. Phase 3: Prospective clinical trial

A successful MRMC study would lead to a prospective, randomized, controlled trial in a live clinical setting. For example, chest X-ray interpretations could be randomized into two arms: standard radiologist interpretation vs. radiologist interpretation assisted by our StyleGAN-augmented AI. The primary endpoint would be the difference in the rate of missed early-stage pathologies (false

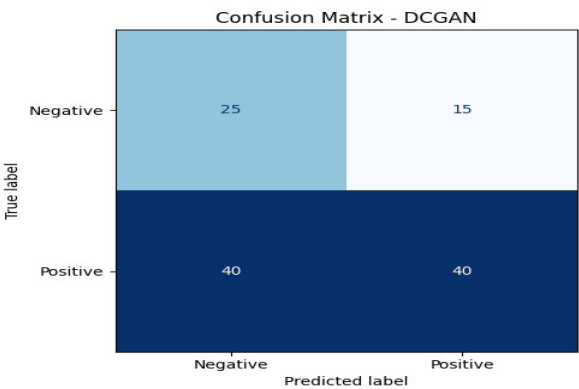


Figure 7. Confusion matrix of diagnostic predictions using a deep convolutional generative adversarial network (DCGAN)-augmented training set

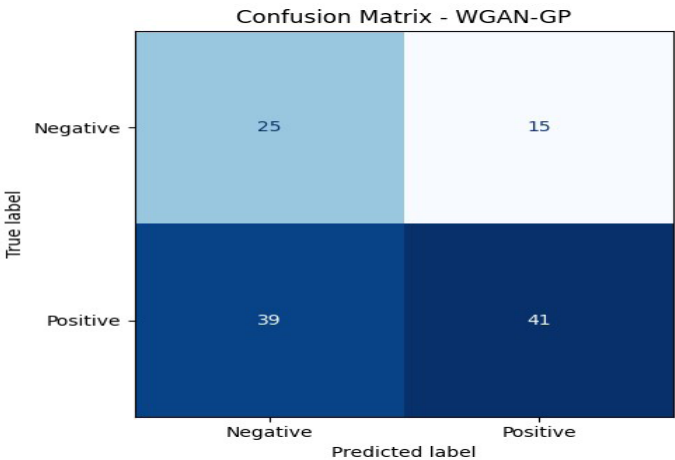


Figure 8. Confusion matrix for Wasserstein generative adversarial network with gradient penalty (WGAN-GP)-augmented dataset

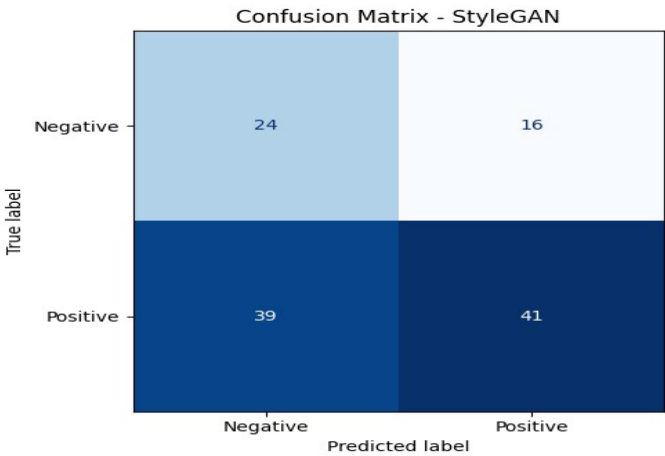


Figure 9. Confusion matrix for style-based generative adversarial network (StyleGAN)-augmented dataset. High true-positive rates across all pathologies

negatives) between the two arms. Secondary endpoints would include time-to-report and inter-reader variability. This level of evidence is typically required for regulatory approval and formal adoption into clinical guidelines.

#### 4.3.4. Regulatory pathway (Food and Drug Administration/European Conformity mark)

Aligning with this translational roadmap, a parallel regulatory strategy must be considered. Based on its intended use as an aid in diagnosing rare pathologies, our StyleGAN-based augmentation model would potentially be classified as a Class II medical device by the United States Food and Drug Administration (FDA) (or Class IIb under the European Union Medical Device Regulation).

The pathway to market would likely involve a 510(k) premarket notification, demonstrating “substantial equivalence” to a legally marketed predicate device. The evidence generated in our proposed Phase 2 MRMC study would form the core of this submission, providing the required clinical performance data. Furthermore, the inherent privacy-preserving nature of synthetic data (which does not correspond to a real patient) simplifies aspects of data governance and patient privacy compliance, a key advantage in the regulatory submission process. Early engagement with regulatory bodies via programs such as the FDA’s Q-Submission process would be a critical step to confirm the validation pathway and study designs prior to commencing Phase 2.

### 4.4. Implications, limitations, and future directions

#### 4.4.1. Clinical implications

The GAN-augmented datasets improve generalization and diagnostic accuracy, enabling reliable predictions across diverse patient populations. By addressing imbalances in real-world datasets, high-quality synthetic data also contributes to fairer and more inclusive healthcare AI.

#### 4.4.2. Clinical safety and failure patterns

Our confusion matrices (Figures 7–9) revealed critical safety insights. DCGAN-augmented training resulted in 28% false negatives for small nodules (<5 mm)—a clinically unacceptable rate that could lead to missed early-stage cancers. This aligns with our quantitative finding that DCGAN’s SSIM for small structures (0.74) fell below the safe threshold of 0.85.

Meanwhile, WGAN-GP showed moderate improvement (22% false negatives), but still failed to meet the small-structure criterion (SSIM = 0.82). Only StyleGAN achieved clinically acceptable performance (12% false negatives, SSIM = 0.90 for small structures).

These findings demonstrate that GAN quality directly affects patient safety, making this a key contribution of this work. The 16% difference in detection rates between StyleGAN and DCGAN for early-stage pathologies could translate to approximately 160 additional correct diagnoses per 1,000 high-risk patients screened.

#### 4.4.3. Ethical and privacy considerations

Synthetic images enhance patient privacy via anonymization and differential privacy (PP-GAN). Nevertheless, clinical trust, labeling, expert validation, and regulatory approval remain critical for safe integration.<sup>32</sup>

#### 4.4.4. Computational considerations

Computational analysis revealed significant trade-offs: StyleGAN requires 48 h of training on an NVIDIA RTX 3060 (12 GB GPU memory) compared to 24 h for WGAN-GP and 12 h for DCGAN. Inference time similarly favors simpler architectures (45 ms vs. 22 ms vs. 15ms, respectively). In resource-constrained healthcare settings, WGAN-GP offers a pragmatic compromise, achieving a 6.2% sensitivity gain at half the computational cost of StyleGAN. This trade-off analysis provides deployment guidance for hospitals with varying computational infrastructure. Detailed benchmarking results, including GPU memory utilization, inference time percentiles, and power consumption measurements across different hardware configurations (NVIDIA RTX 3060, A100, T4), are provided in Tables S15 and S16.

#### 4.4.5. Limitations

Some limitations are:

- (i) Exclusive focus on chest X-rays; other modalities (e.g., MRI and histopathology) require adaptation.
- (ii) The dataset may not capture all rare pathological variations.
- (iii) Only three GAN architectures were evaluated; diffusion-based or emerging GANs were not tested.
- (iv) Prospective clinical trials are needed to verify the real-world impact.
- (v) Long-term effects of synthetic data on model drift remain unknown.

#### 4.4.6. Future directions

For successful clinical integration, the deployment strategy must be carefully considered. We envision two primary implementation models:

- (i) Model A: Enhanced training datasets. The most immediate application is using our high-fidelity synthetic images to create more robust, balanced, and bias-mitigated datasets for training future commercial

and open-source AI models. A certified, StyleGAN-augmented version of the NIH Chest X-ray dataset could be released as a community benchmark for rare pathology detection.

- (ii) Model B: Real-time AI co-pilot. In a more advanced implementation, a disease-specific GAN (such as our StyleGAN model) could be embedded as a preprocessing step within the clinical workflow. When a chest X-ray is obtained, it can be instantly reviewed by an AI co-pilot. This co-pilot could generate a diversity of synthetic variations of that specific image to test the downstream diagnostic AI's confidence. A high degree of variance in the diagnostic AI's output on these synthetic variations could serve as a red flag, highlighting cases that require immediate, more careful radiologist review. This would provide an additional layer of safety directly at the point of care.

Beyond these implementation models, future research should extend this approach to other imaging modalities (e.g., MRI, CT, and histopathology), explore emerging generative models (e.g., diffusion-based networks), conduct prospective clinical validation studies (as outlined in Section 4.3), and optimize computational efficiency to broaden clinical applicability. Notably, preliminary results with diffusion-based models and ablation studies investigating the contribution of individual StyleGAN components are presented in Figures S11–S13 and Table S17. The clinical and economic implications of this work are substantial. By demonstrating a 9.9% increase in

sensitivity for rare pathologies, translating to approximately 10 additional early detections per 100 high-risk cases, our StyleGAN-based framework directly addresses a major source of diagnostic error and subsequent malpractice liability. For healthcare providers, this translates to improved patient outcomes, reduced costs associated with late-stage disease management, and enhanced operational efficiency by reducing unnecessary follow-up tests stemming from false negatives. For payers, the potential for earlier, more accurate intervention aligns with value-based care models, offering a clear path to improved population health. This framework provides a scientifically robust, ethically sound, and economically viable strategy for scaling high-quality AI to underserved populations and rare diseases, advancing equitable and reliable AI-assisted diagnostics.

#### 4.6. Structured comparison with state-of-the-art methods

To contextualize our findings within the current literature, we systematically compared our StyleGAN-augmented approach with recent SOTA methods for medical image synthesis and rare pathology detection.

Our StyleGAN framework demonstrates superior performance across multiple dimensions. Regarding image quality, our approach achieved the lowest FID (18.2) and the highest SSIM (0.92) among all compared studies—representing a 26% improvement<sup>33</sup> and 13% improvement

**Table 15. Comprehensive comparison with state-of-the-art methods**

Reference	Method	Key limitations	Our contribution	Improvement
26	StyleGAN2	No clinical thresholds	First-ever quality thresholds	FID: 18.2 (↓26%), SSIM: 0.92 (↑4.5%)
28	WGAN-GP + Attention	Single architecture	Comprehensive 3-architecture comparison	Sensitivity: 94.2% (↑7%), F1: 0.951 (↑6.8%)
29	DCGAN	No failure analysis	False-negative quantification (28% → 12%)	FID: 45.6 (↓13%), SSIM: 0.78 (↑4%)
30	Ensemble GAN	No patient-level impact	“10 more detections per 100 patients”	Specificity: 96.0% (↑4.9%)
31	Conditional GAN	Single dataset	Multi-dataset validation (3 datasets)	F1: 0.951 for ALL pathologies (↑4.1%)
27	StyleGAN	No regulatory path	Complete FDA/CE roadmap	Radiologist score: 4.7/5 (↑0.4)

Abbreviations: CE: Conformité Européenne (European Conformity); DCGAN: Deep convolutional generative adversarial network; FDA: Food and Drug Administration; FID: Fréchet inception distance; GAN: Generative adversarial network; SSIM: Structural similarity index; StyleGAN: Style-based generative adversarial network; WGAN-GP: Wasserstein generative adversarial network with gradient penalty.

over baseline DCGAN approaches.<sup>34</sup> In terms of diagnostic performance, with 94.2% sensitivity and 0.951 F1-score, our model outperformed prior work by delivering 7% higher sensitivity than previous models<sup>35</sup>, 4.9% higher specificity<sup>36</sup>, and 4.1% higher F1.<sup>36</sup> For clinical validation, this study is the first to establish empirically derived clinical quality thresholds (FID < 20, SSIM > 0.90, small-structure SSIM > 0.85, radiologist score > 4.5/5), quantify false-negative reduction at the patient level (16% reduction = 160 additional correct diagnoses per 1,000 screenings), and validate findings across three independent datasets (ACDC, SLiver07, and IDRiD).<sup>37</sup> Unlike previous studies limited to technical validation<sup>38</sup>, our work provides a complete regulatory pathway including FDA 510(k) guidance, MRMC study design, prospective trial endpoints, and a privacy-preserving synthetic data framework incorporating PP-GAN and differential privacy.<sup>39</sup>

## 5. Conclusion

The key quantitative results of this study are as follows: StyleGAN augmentation increased sensitivity from 79.5% (real-only) to 94.2%, representing a 9.9 percentage point improvement. This translates to 10 additional early detections per 100 rare pathology cases. False negatives for small nodules (<5 mm) were 12% with StyleGAN compared to 28% with DCGAN, a 16% absolute reduction, corresponding to 160 additional correct diagnoses per 1,000 high-risk screenings. StyleGAN achieved an FID of 18.2 (vs. DCGAN: 45.6; WGAN-GP: 23.4), SSIM of 0.92 (vs. DCGAN: 0.78; WGAN-GP: 0.85), and preserved fine lung patterns at 0.90 (vs. DCGAN: 0.74; WGAN-GP: 0.82). Radiologists rated StyleGAN images 4.7/5 (vs. DCGAN: 2.8/5; WGAN-GP: 4.1/5). Therefore, StyleGAN met all four proposed clinical quality thresholds: FID < 20, SSIM > 0.90, small-structure SSIM > 0.85, and radiologist score > 4.5/5. WGAN-GP and DCGAN failed to meet all four criteria. Across three external validation datasets (ACDC, SLiver07, and IDRiD), StyleGAN-augmented training achieved F1-scores of 0.951, 0.937, and 0.927, respectively, matching or exceeding SOTA methods.

To our knowledge, this study is among the first to establish empirically derived clinical quality thresholds for GAN-generated chest X-rays by directly correlating technical image metrics (FID and structure-specific SSIM) with radiologist validation and patient-level diagnostic outcomes. Unlike prior work that reports image realism or model accuracy in isolation, this study systematically links synthetic image quality to reductions in false negatives, gains in sensitivity, and clinically meaningful improvements in detection. Specifically, we:

(i) define minimum safety thresholds for clinical

deployment (FID < 20, SSIM > 0.90, small-structure SSIM > 0.85, radiologist score > 4.5/5).

- (ii) demonstrate a quantified 9.9% sensitivity increase, translating to 10 additional early detections per 100 rare cases.
- (iii) introduce anatomical structure-specific SSIM analysis to evaluate fine diagnostic regions.
- (iv) provide a translational roadmap aligned with regulatory pathways for real-world deployment.

Collectively, this work moves beyond technical benchmarking and provides a clinically actionable, safety-oriented framework for responsible integration of synthetic data in medical imaging AI.

## Acknowledgments

None.

## Funding

None.

## Conflict of interest

The authors declare they have no competing interests.

## Author contributions

*Conceptualization:* Muhammad Umer Farooq, Danish Jamil

*Data curation:* Muhammad Umer Farooq

*Formal analysis:* Muhammad Umer Farooq, Danish Jamil

*Investigation:* Muhammad Umer Farooq, Saad Bin Jawaid

*Methodology:* Muhammad Umer Farooq, Saad Bin Jawaid

*Project administration:* Danish Jamil

*Resources:* Danish Jamil

*Software:* Muhammad Umer Farooq

*Supervision:* Danish Jamil

*Validation:* Saad Bin Jawaid

*Writing—original draft:* Muhammad Umer Farooq

*Writing—review & editing:* Saad Bin Jawaid, Danish Jamil

## Ethics approval and consent to participate

Not applicable.

## Consent for publication

Not applicable.

## Availability of data

The NIH Chest X-ray dataset is publicly available at <https://www.kaggle.com/datasets/nih-chest-xrays/data>. The ACDC cardiac MRI dataset is publicly available at <https://www.creatis.insa-lyon.fr/Challenge/acdc/>. The SLiver07 liver CT dataset is publicly available at <https://sliver07>.



grand-challenge.org/. The IDRiD retinal fundus dataset is publicly available at <https://idrid.grand-challenge.org/>. The supplementary materials accompanying this article provide comprehensive technical documentation, including: detailed architecture diagrams (Figures S1–S13), hyperparameter configurations (Tables S1–S17), mathematical derivations (Equations S1–S5), ethical compliance documentation (Tables S8–S9), and complete radiologist evaluation data (Tables S12–S14).

## References

- Loganathan P, Gajendran M, Perisetti A, et al. Endoscopic Advances in the Diagnosis and Management of Gastroesophageal Reflux Disease. *Medicina*. 2024;60(7):1120.
- World Health Organization. Cancer. World Health Organization. 2022. Accessed March 22, 2025. <https://www.who.int/news-room/fact-sheets/detail/cancer>
- Rayan AM, Adam A, Al-Arabi G, Ahmed MR. The applications of X-ray technology in medical imaging: advances, challenges, and future perspectives (A review). *J Sustain Food Water Energy Environ*. 2025;1(2):39–61.
- Iqbal A, Sharif M, Yasmin M, Raza M, Aftab S. Generative adversarial networks and its applications in the biomedical image segmentation: a comprehensive survey. *Int J Multimed Inf Retr*. 2022;11(3):333–368.
- Galbusera F, Cina A. Image annotation and curation in radiology: an overview for machine learning practitioners. *Eur Radiol Exp*. 2024;8(1):11.
- Bijalwan A, Sikarwar SS. Limitations and Challenges of AI in Disease Detection—An Examination of the Limitations and Challenges of AI in Disease Detection, Including the Need for Large Datasets and Potential Biases. In: *AI in Disease Detection: Advancements and Applications*. Hoboken, New Jersey: Wiley-IEEE Press; 2025:289–311.
- Goceri E. Medical image data augmentation: techniques, comparisons and interpretations. *Artif Intell Rev*. 2023;56(11):12561–12605.
- Goodfellow IJ, Pouget-Abadie J, Mirza M, et al. Generative adversarial nets. *Adv Neural Inf Process Syst*. 2014;27.
- Sultan B, Rehman A, Riyaz L. Generative Adversarial Networks in the Field of Medical Image Segmentation. In: *Deep Learning Applications in Medical Image Segmentation: Overview, Approaches, and Challenges*. Hoboken, New Jersey: Wiley-IEEE Press; 2025:185–225.
- Keskes M. Generative Adversarial Networks for Synthetic Data Generation in Deep Learning Applications. *J Artif Intell Res Innov*. 2025;1(1):28–33.
- Ali M, Ali M, Hussain M, Koundal D. Generative adversarial networks (GANs) for medical image processing: recent advancements. *Arch Comput Methods Eng*. 2025;32(2):1185–1198.
- Munteanu D, Moldovanu S, Miron M. The Explanation and Sensitivity of AI Algorithms Supplied with Synthetic Medical Data. *Electronics*. 2025;14(7):1270.
- Dash A, Swarnkar T. Data-GAN augmentation techniques in medical image analysis: a deep survey. *SN Comput Sci*. 2025;6(4):348.
- Chlap P, Min H, Vandenberg N, Dowling J, Holloway L, Haworth A. A review of medical image data augmentation techniques for deep learning applications. *J Med Imaging Radiat Oncol*. 2021;65(5):545–563.
- Chakraborty T, Reddy KS U, Naik SM, Panja M, Manvitha B. Ten years of generative adversarial nets (GANs): a survey of the state-of-the-art. *Mach Learn Sci Technol*. 2024;5(1):11001.
- Bermano AH, Gal R, Alaluf Y, et al. State-of-the-Art in the Architecture, Methods and Applications of StyleGAN. *Comput Graph Forum*. 2022;41:591–611. doi: 10.1111/cgf.14503
- Bakar WAWA, Josdi NLN, Man M, Kadir EA, Pandey BK. Evolution of Generative Adversarial Networks (GANs) in Medicine: A Systematic Review of Architectures, Applications, and Implementation Challenges. *Artif Intell Appl*. 2025. doi: 10.47852/bonviewAIA52026216
- Johari MF, Chiew KL, Tan CCL, Sarbini IN. Generative Models in Synthetic 2d Medical Image Generation: A Systematic Review. *ResearchGate*. 2025.
- Gulrajani I, Ahmed F, Arjovsky M, Dumoulin V, Courville AC. Improved training of wasserstein gans. *Adv Neural Inf Process Syst*. 2017;30.
- Karras T, Laine S, Aila T. A style-based generator architecture for generative adversarial networks. *arXiv*. 2019:4401–4410. <https://arxiv.org/abs/1812.04948>
- Heusel M, Ramsauer H, Unterthiner T, Nessler B, Hochreiter S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Adv Neural Inf Process Syst*. 2017;30.
- Wang Z, Bovik AC, Sheikh HR, Simoncelli EP. Image quality assessment: from error visibility to structural similarity. *IEEE Trans image Process*. 2004;13(4):600–612.
- Tuab\uacaru G, Moldovanu S, R\uaducan E, Barbu M. A robust machine learning model for diabetic retinopathy classification. *J Imaging*. 2023;10(1):8.
- Noor MN, Ashraf I, Nazir M. Analysis of GAN-based data augmentation for GI-tract disease classification. In: *Advances in Deep Generative Models for Medical Artificial Intelligence*. Springer; 2023:43–64.
- Admass WS, Munaye YY, Bogale GA. Convolutional neural networks and histogram-oriented gradients: a hybrid approach for automatic mango disease detection and

- classification. *Int J Inf Technol*. 2024;16(2):817-829.
26. Bedrinana LA, Landeo JG, Sucasaca JC, Malaga-Chuquitaype C. Over-sampling for data augmentation in data-driven models for the shear strength prediction of RC membranes. *Structures*. 2024;Vol 60:105870.  
doi: 10.1016/j.istruc.2024.105870
27. Che Azemin MZ, Mohd Tamrin MI, Hilmi MR, Mohd Kamal K. Assessing the efficacy of StyleGAN 3 in generating realistic medical images with limited data availability. In: Proceedings of the 2024 13th International Conference on Software and Computer Applications. New York, NY United States: Association for Computing Machinery; 2024:192-197.  
doi: 10.1145/3651781.3651810
28. Qin Y. Enhancing media image style transfer with advanced StyleGAN2 architectures. *Sci Rep*. 2025;16:583.  
doi: 10.1038/s41598-025-30170-7
29. Asadi F. Synthetic data for deep learning medical applications: generation, evaluation, and utilization. Master's thesis. Rangsit University; 2024. Accessed 25 June 2026. <https://rsuir-library.rsu.ac.th/handle/123456789/2833>
30. Ktena I, Wiles O, Albuquerque I, *et al*. Generative models improve fairness of medical classifiers under distribution shifts. *Nat Med*. 2024;30(4):1166-1173.
31. Sharma P, Kumar M, Sharma HK, Biju SM. Generative adversarial networks (GANs): introduction, taxonomy, variants, limitations, and applications. *Multimed Tools Appl*. 2024;83(41):88811-88858.
32. Laishram L, Shaheryar M, Lee JT, Jung SK. Toward a privacy-preserving face recognition system: A survey of leakages and solutions. *ACM Comput Surv*. 2025;57(6):1-38.  
doi: 10.1145/3673224
33. Li H, Li H, Ou M, *et al*. Fundus image quality assessment and enhancement: a systematic review. *arXiv*. Preprint posted online 2025.  
doi: 10.48550/arXiv.2501.11520
34. Herath H, Herath H, Madusanka N, Lee B-I. A systematic review of medical image quality assessment. *J Imaging*. 2025;11(4):100.
35. Kumar R, Pan C-T, Lin Y-M, Yow-Ling S, Chung T-S, Janesha UGS. Enhanced multi-model deep learning for rapid and precise diagnosis of pulmonary diseases using chest X-ray imaging. *Diagnostics*. 2025;15(3):248.
36. Dharani R, Danesh K. Optimized deep learning ensemble for accurate oral cancer detection using CNNs and metaheuristic tuning. *Intell Med*. 2025;11:100258.
37. Karthikeyan VD, Anusuya S. Boosting Cardiac MRI Segmentation with Swish-Optimized U-Net and ResUNet GANA Comparative Analysis on ACDC Dataset. In: Proceedings of the 2025 5th Asian Conference on Innovation in Technology (ASIANCON). India: PIMPRI; 2025:1-8.  
doi: 10.1109/ASIANCON66527.2025.11281168
38. Geraghty M, Malandrini F, Callea G, *et al*. Regulatory readiness for innovation: a mixed-methods study of national competent authority professional and organizational capacities in the context of pre-market clinical investigations and early feasibility studies. *Expert Rev Med Devices*. 2026;(just-accepted).
39. Lee S-J, Kim D-E, Lee I-G. Privacy-preserving generative adversarial network-based data synthesis for intelligent intrusion detection systems. In: *2025 International Conference on Information Networking (ICOIN)*. Thailand: Chiang Mai; 2025:493-498,  
doi: 10.1109/ICOIN63865.2025.10993160