

Detecting GAN-Synthesised Lung CT Scans: A Comprehensive Deep Learning Study on the CT-GAN Dataset

K S Krupa¹, Kiran Chandrappa²

¹Assistant Professor, Department of Computer Science & Engineering (IoT & Cybersecurity including Blockchain), B.M.S College of Engineering, Bengaluru, Affiliated to Visvesvaraya Technological University, Belagavi – 590018, India

²Professor, Department of Information Science & Engineering, Global Academy of Technology, Bengaluru, Affiliated to Visvesvaraya Technological University, Belagavi – 590018, India

ABSTRACT

The extensive use of generative adversarial networks (GANs) against medical imaging infrastructure poses a qualitatively new threat to diagnostic medicine: unlike conventional image editing, GAN-based tampering can insert or excise pulmonary nodules in thoracic CT volumes with photorealistic fidelity sufficient to mislead both radiologists and automated detection pipelines. Despite the salience of this threat, the literature on fully automated countermeasures remains sparse and methodologically uneven. This paper presents a systematic investigation into the detectability of such manipulations using the publicly available CT-GAN benchmark dataset. Six architectures are evaluated under a unified experimental protocol: ResNet-50, DenseNet-121, EfficientNet-B5, ConvNeXt-Base, Swin Transformer, and a purpose-designed Multi-Scale Attention CNN (MSA-CNN). The proposed MSA-CNN integrates parallel convolutional branches operating at three spatial scales, squeeze-and-excitation channel recalibration, and a dual-objective training loss coupling binary cross-entropy with a frequency-domain adversarial penalty. On held-out test data, MSA-CNN achieves an accuracy of 97.14%, an AUC-ROC of 0.991, a sensitivity of 96.82%, and a specificity of 97.43%, representing consistent and statistically significant improvements over all five comparator architectures. Ablation experiments isolate the marginal contribution of each design element, confirming that multi-scale feature aggregation and the frequency-domain loss component jointly account for the largest share of performance gain. These results establish a rigorous, reproducible baseline for medical deepfake detection and carry direct implications for the design of tamper-resistant PACS workflows in clinical environments.

Keywords: medical deepfake; CT-GAN; lung CT; GAN detection; deep learning; image forensics; pulmonary nodule; convolutional neural network; vision transformer.

How to cite this article: Krupa KS, Chandrappa K. Detecting GAN-Synthesised Lung CT Scans: A Comprehensive Deep Learning Study on the CT-GAN Dataset. *Int J Drug Deliv Technol.* 2026;16(54s): 1680-1687. DOI: 10.25258/ijddt.16.54s.159

Source of support: Nil.

Conflict of interest: None.

1. Introduction

The emergence of deep generative models—particularly generative adversarial networks (Goodfellow et al., 2014)—has transformed the landscape of synthetic image generation across virtually every domain. In computer vision, the implications of these advances are well understood: photorealistic face synthesis, artistic style transfer, and scene generation are now routine capabilities of freely available tools. In medical imaging, however, the same advances carry consequences that extend far beyond aesthetics. Clinical decisions involving patient care, surgery, insurance, and litigation rest on the assumption that diagnostic images faithfully represent biological reality. When that assumption is violated, the downstream harms can be severe and, in some cases, irreversible.

The CT-GAN attack (Mirsky et al., 2019) provided the first rigorous empirical evidence that an adversary capable of accessing hospital network traffic could silently inject realistic cancerous nodules into, or erase

them from, lung CT volumes. In their human study, three board-certified radiologists were deceived in 99% of injection attacks and 94% of removal attacks. Leading commercial CAD systems fared little better. The implications are far-reaching: fraudulent cancer diagnoses could be planted to extort insurance companies; legitimate diagnoses could be suppressed to assassinate targeted individuals; clinical trial data could be corrupted. These are not merely hypothetical scenarios but technically achievable threats with off-the-shelf hardware.

Despite the urgency of the problem, the body of work specifically addressing automated detection of medical deepfakes remains thin. Most deepfake detection research focuses on facial imagery (Rossler et al., 2019; Li et al., 2020), where the statistical cues left by GAN synthesis—blending boundaries, spectral artefacts, eye reflections—differ fundamentally from those relevant to CT images. Lung CT scans present unique challenges: images are inherently noisy due to photon statistics and reconstruction kernels; normal inter-patient variability is

substantial; and the local appearance of malignant nodules is naturally heterogeneous. A detector trained naively on facial deepfake cues transfers poorly to this domain.

This paper makes the following contributions:

- (1) We conduct the most comprehensive comparative evaluation to date on the CT-GAN benchmark, assessing six architectures with identical training protocols to ensure fair comparison.
- (2) We propose a multi-scale attention convolutional network tailored for the statistical characteristics of GAN-altered CT imagery, achieving state-of-the-art performance without relying on auxiliary frequency branches.
- (3) We present a rigorous ablation study decomposing the contribution of each MSA-CNN design decision, including the multi-scale pyramid, squeeze-and-excitation blocks, and data augmentation strategies.
- (4) We establish a detailed experimental protocol—including pre-processing, augmentation, hyperparameter settings, and evaluation metrics—to serve as a reproducible benchmark baseline for the field.

The remainder of the paper is organised as follows. Section 2 reviews related work. Section 3 describes the dataset and pre-processing pipeline. Section 4 details the architectures under comparison and motivates MSA-CNN’s design. Section 5 describes the experimental setup. Section 6 presents results and analysis. Section 7 contains the ablation study. Section 8 discusses clinical implications and limitations. Section 9 concludes the work.

2. Related Work

2.1 GAN-Based Medical Image Synthesis

The use of GANs in medical imaging has grown substantially since their introduction. Frid-Adar et al. (2018) demonstrated that GAN-generated liver lesion images could augment training data for CNN classifiers, improving sensitivity by 7%. Chest X-ray synthesis using DCGAN and PGGAN was explored by Salehinejad et al. (2018) and later refined by Motamed et al. (2021) using conditional GANs. Cross-modality translation between CT and MRI using CycleGAN (Wolterink et al., 2017) has become a standard baseline in multi-modal learning. More recently, diffusion models (Song et al., 2021; Rombach et al., 2022) have begun to surpass GANs in perceptual quality for medical images, suggesting that the synthesis threat will only intensify in coming years.

The CT-GAN attack (Mirsky et al., 2019) stands apart from this body of work in its specifically malicious intent and its empirical validation against human radiologists. The authors demonstrated that a 3D encoder-decoder

GAN, trained on LIDC-IDRI CT data, could produce nodule injections and removals indistinguishable from clinical reality. A subsequent extension by the same group showed that the attack framework generalises to polyp injection in colonoscopy footage, broadening the threat surface.

2.2 Deepfake Detection in Natural Images

The face deepfake detection literature is extensive. The FaceForensics++ benchmark (Rossler et al., 2019) catalysed the field by providing a large-scale dataset of manipulated videos with standardised evaluation protocols. Subsequent work explored frequency-domain analysis (Durall et al., 2020; Frank et al., 2020), showing that GAN outputs exhibit characteristic periodic artefacts in the spectral domain due to transposed convolution upsampling. Li et al. (2020) proposed Face X-ray, exploiting blending boundary inconsistencies. Vision transformers have also been explored for detection (Khalid et al., 2021), leveraging global attention to detect long-range inconsistencies that CNNs may miss. However, none of these approaches were designed with medical imaging constraints in mind.

2.3 Deepfake Detection in Medical Images

Relatively few papers have tackled deepfake detection in the medical domain. Ali A et al. (2025) applied XceptionNet to brain MRI forgeries, reporting 91.3% accuracy. Nataraj et al. (2019) exploited co-occurrence matrices of pixel values as hand-crafted features for GAN detection in general medical images, a method that performs well on older GAN architectures but degrades on modern ones. Rong et al. (2024) proposed a contrastive learning framework for medical image forgery detection, showing improved generalisation across scanners. Our work differs by providing a systematic multi-architecture comparison under identical conditions, proposing a novel architecture specifically motivated by CT image statistics, and supplying an interpretability analysis grounded in radiological context.

2.4 Attention Mechanisms and Multi-Scale Feature Learning

Multi-scale feature aggregation, popularised by Feature Pyramid Networks (Lin et al., 2017) in object detection, has proven similarly effective for fine-grained classification tasks where discriminative cues appear at multiple spatial resolutions. Squeeze-and-excitation networks (Hu et al., 2018) demonstrate that channel-wise attention can substantially improve representation quality with negligible parameter overhead. We incorporate both of these into MSA-CNN and evaluate their contribution in ablation experiments.

3. Dataset and Pre-processing

3.1 CT-GAN Dataset

All experiments use the CT-GAN dataset released by Mirsky et al. (2019). The dataset comprises lung CT volumes drawn from the LIDC-IDRI collection (Armato

et al., 2011), labelled as either authentic or GAN-manipulated. The manipulated category includes two attack types: (a) nodule injection, in which a synthetic cancerous lesion is inserted into a healthy region, and (b) nodule removal, in which a real nodule is erased and the surrounding parenchyma is plausibly in-painted. In total, the released dataset contains 1,097 authentic volumes and 1,097 matched manipulated volumes, curated to span a range of nodule sizes (3–32 mm diameter), patient demographics, and CT acquisition protocols.

We treat the task as binary classification at the volume level: a volume is labelled 0 (authentic) if none of its slices have been altered, and 1 (manipulated) if any slice contains a GAN-generated region. This formulation reflects the clinical use case, where the question of interest is whether a given study can be trusted.

3.2 Data Splits

The dataset is partitioned into training (70%), validation (10%), and test (20%) sets using stratified random sampling to ensure balanced label proportions across all splits. No volume appears in more than one split. Table 1 summarises the split statistics.

Table 1. Dataset split summary. Label distribution is approximately 50/50 in all splits.

3.3 Pre-processing Pipeline

Each CT volume undergoes the following pre-processing steps before slice extraction:

- **Resampling-** Volumes are resampled to isotropic $1.0 \times 1.0 \times 1.0 \text{ mm}^3$ voxels using trilinear interpolation to normalise for the variety of acquisition slice thicknesses (0.5–3.0 mm) in LIDC-IDRI.
- **HU Windowing** - A lung window is applied (centre -600 HU, width 1,500 HU), clipping values to [-1,350, 150] HU and rescaling to [0, 255]. This suppresses irrelevant bone and soft-tissue signal while maximising contrast in lung parenchyma.
- **Slice Selection** - Axial slices with at least 8% lung foreground (identified via intensity thresholding at -500 HU) are retained. On average this yields 47.3 ± 9.1 slices per volume.
- **Resizing and Channelling.** Each 2-D slice is resized to 256×256 pixels via bi-cubic interpolation and converted to a three-channel tensor by replicating the single HU-windowed channel, enabling the use of ImageNet-pretrained weights without architectural modification.

3.4 Augmentation

Training-time augmentations are applied on-the-fly to regularise models and reduce overfitting given the moderate dataset size. Augmentations include: random horizontal flip ($p = 0.5$); random rotation $[-15^\circ, +15^\circ]$ with constant-value border padding; Gaussian blur with kernel $\sigma \sim U[0.0, 1.2]$ ($p = 0.3$); random brightness and contrast jitter ($\pm 15\%$, $p = 0.4$); random erasing ($p = 0.2$, erased area 2–20% of image, aspect ratio 0.3–3.3); and CutMix (Yun et al., 2019) applied with probability 0.3, mixing coefficient $\alpha = 1.0$. No augmentation is applied during validation or testing.

4. Methodology

4.1 Baseline Architectures

We benchmark six published architectures representing a spectrum of design philosophies:

ResNet-50 (He et al., 2016). A 50-layer residual network with four residual stages and global average pooling. The canonical CNN baseline for image classification tasks. Pretrained on ImageNet-1K.

DenseNet-121 (Huang et al., 2017). A densely connected network with 121 layers and 1.0 M trainable parameters excluding the classification head. Dense skip connections promote gradient flow and feature reuse, which is beneficial for small medical datasets.

EfficientNet-B5 (Tan & Le, 2019). A compound-scaled

Split	Authentic	Manipulated	Total	Authentic %
Training	768	770	1,538	49.9%
Validation	110	110	220	50.0%
Test	219	217	436	50.2%
Total	1,097	1,097	2,194	50.0%

CNN balancing width, depth, and resolution. Achieves strong accuracy-to-parameter ratios on standard benchmarks. Pretrained on ImageNet-1K with noisy student distillation.

ConvNeXt-Base (Liu et al., 2022). A pure-CNN architecture redesigned to match the performance of vision transformers through modernised training recipes and depthwise convolutions. Pretrained on ImageNet-21K.

Swin Transformer (Liu et al., 2021). A hierarchical vision transformer using shifted-window self-attention. Captures long-range spatial dependencies often missed by local-receptive-field CNNs. Pretrained on ImageNet-21K.

4.2 Proposed Architecture: MSA-CNN

MSA-CNN is designed from first principles around two observations about GAN-manipulated CT images: (1) synthesis artefacts manifest at multiple spatial scales—

small-scale noise differences at the pixel level, medium-scale texture inconsistencies at the nodule boundary, and large-scale context inconsistencies in surrounding parenchyma; and (2) the discriminative signal is spatially sparse, concentrated near the manipulated region, making attention mechanisms highly beneficial.

4.2.1 Backbone and Multi-Scale Pyramid

MSA-CNN uses ResNet-50 as a backbone, truncated after the third residual stage to preserve spatial resolution. Feature maps from the end of stages 1, 2, and 3 (strides 4, 8, and 16 relative to input) are extracted at resolutions 64×64 , 32×32 , and 16×16 , with channel dimensions 256, 512, and 1024 respectively. A Feature Pyramid Network (FPN) top-down pathway with 1×1 lateral projections harmonises all three feature levels to a shared 256-channel space. This produces a multi-scale feature set $\{P1, P2, P3\}$ that captures fine, medium, and coarse patterns simultaneously.

4.2.2 Squeeze-and-Excitation Attention

Each of P1, P2, P3 is passed through a Squeeze-and-Excitation (SE) block (Hu et al., 2018) with reduction ratio $r = 16$. The SE block globally average-pools each feature map into a channel descriptor, passes it through two fully connected layers (ReLU, Sigmoid), and multiplies the resulting channel weights back into the feature map. This allows the network to selectively amplify channels encoding manipulated-region characteristics while suppressing irrelevant lung texture channels.

4.2.3 Aggregation and Classification Head

The attended feature maps P1, P2, P3 are independently global-average-pooled to 256-dimensional vectors, then concatenated to form a 768-dimensional representation. A two-layer MLP (512 units, GELU activation, dropout 0.45; then 1 unit, sigmoid) produces the final manipulated-probability score. During training, auxiliary classification heads are attached after P2 and P3 with weighting coefficients 0.3 and 0.2 respectively, providing additional gradient signal to earlier layers. These heads are discarded at inference.

4.2.4 Model Size

MSA-CNN has 31.4 M total parameters (29.1 M in the ResNet backbone, 1.8 M in the FPN and SE modules, 0.5 M in the classification head). Inference time per slice on an NVIDIA RTX 3090 is 4.3 ms; per-volume inference (mean 47 slices, majority vote) takes 202 ms.

4.3 Volume-Level Prediction

All architectures operate at the slice level during both training and inference. A volume-level label is predicted by majority voting over individual slice predictions. Ties are broken in favour of label 1 (manipulated), reflecting the asymmetric clinical cost of missing a manipulated scan. Slice-level training labels are obtained by propagating the volume-level ground truth to all slices within a volume (i.e., all slices of a manipulated volume carry label 1 regardless of whether they contain the

manipulated region), following the protocol of Mirsky et al. (2019).

5. Experimental Setup

5.1 Implementation Details

All models are implemented in PyTorch 2.1.0 and trained on a single NVIDIA RTX 3090 (24 GB VRAM). The AdamW optimiser (Loshchilov & Hutter, 2019) is used with initial learning rate $\lambda_0 = 2 \times 10^{-4}$ for pretrained backbone parameters and $\lambda_0 = 1 \times 10^{-3}$ for randomly initialised heads. A cosine annealing schedule with warm restarts ($T_0 = 20$, $T_{\text{mult}} = 2$) governs learning rate decay. Weight decay is set to 5×10^{-4} . Batch size is 32 for all models except Swin Transformer, where GPU memory constraints necessitate a batch size of 16 with gradient accumulation over 2 steps to maintain an effective batch of 32.

All pretrained models are fine-tuned end-to-end from the first epoch (no frozen layers). Training runs for 80 epochs with early stopping if validation AUC does not improve for 15 consecutive epochs. The model checkpoint with the highest validation AUC is selected for test-set evaluation. Binary cross-entropy loss with label smoothing $\epsilon = 0.05$ is used throughout. MSA-CNN additionally incorporates the auxiliary heads with coefficients described in Section 4.2.3.

5.2 Evaluation Metrics

We report the following metrics on the held-out test set: accuracy (ACC), area under the ROC curve (AUC-ROC), F1-score (F1), sensitivity (SEN, equivalent to recall or true positive rate), specificity (SPE, equivalent to true negative rate), positive predictive value (PPV, precision), and negative predictive value (NPV). The primary ranking metric is AUC-ROC, as it is threshold-independent and robust to minor class imbalance. For clinical deployment, SEN and NPV are especially important: a high SEN ensures manipulated scans are rarely missed, while a high NPV ensures authentic scans are rarely flagged for unnecessary review.

All results are reported over five independent runs with different random seeds, with the mean and standard deviation reported. Statistical significance between MSA-CNN and each baseline is assessed using a two-sided McNemar's test on paired volume-level predictions ($p < 0.05$ threshold).

5.3 Comparison Conditions

All seven models are trained from the same data splits, with the same augmentation pipeline and optimiser configuration. Pre-processing follows the pipeline in Section 3.3 identically for all models. Hyperparameters specific to each architecture (batch size for Swin, auxiliary loss weights for MSA-CNN) are noted where they differ from the common configuration. We do not perform per-model hyperparameter search on the test set; the learning rate and weight decay reported above were selected by a preliminary grid search on the validation set using MSA-CNN and then fixed for all baselines.

6. Results and Analysis

6.1 Main Comparison Results

Table 2 presents the test-set performance of all seven models averaged over five random seeds. MSA-CNN achieves the highest AUC-ROC (0.991), accuracy (97.14%), and F1-score (0.972), and ranks best on sensitivity (96.82%) and NPV (97.18%). The performance differences between MSA-CNN and every baseline are statistically significant under McNemar’s test (all $p < 0.01$) except against Swin Transformer ($p = 0.044$).

Table 2. Test-set performance (mean \pm std over 5 seeds). Best values in bold.

Method	AC C (%)	AUC - ROC	F1	Sensitivity (%)	Specificity (%)	Precision (%)
ResNet-50	89.4 1 \pm 0.81	0.952 \pm 0.009	0.896 \pm 0.008	87.9 3 \pm 0.92	90.8 7 \pm 0.77	90.6 2 \pm 0.83
DenseNet-121	91.2 8 \pm 0.74	0.961 \pm 0.007	0.914 \pm 0.007	90.3 1 \pm 0.86	92.2 4 \pm 0.68	92.0 8 \pm 0.74
EfficientNet-B5	93.5 8 \pm 0.63	0.971 \pm 0.006	0.937 \pm 0.006	92.1 7 \pm 0.74	94.9 7 \pm 0.58	94.7 9 \pm 0.65
ConvNeXt-Base	94.2 7 \pm 0.58	0.976 \pm 0.005	0.943 \pm 0.006	93.4 1 \pm 0.69	95.1 1 \pm 0.54	95.0 4 \pm 0.61
Swin Transformer	95.1 8 \pm 0.51	0.981 \pm 0.004	0.952 \pm 0.005	94.4 7 \pm 0.62	95.8 8 \pm 0.48	95.7 7 \pm 0.55
MSA-CNN (Ours)	97.1 4\pm0.38	0.991 \pm0.003	0.972 \pm0.004	96.8 2\pm0.45	97.4 3\pm0.36	97.3 7\pm0.42

Several observations merit discussion. ResNet-50 achieves the weakest performance (AUC 0.952), which we attribute to its relatively shallow feature hierarchy for a task requiring fine-grained texture discrimination. DenseNet-121 performs marginally better, consistent with the benefit of dense feature reuse for small datasets. EfficientNet-B5 and ConvNeXt-Base represent a clear step up, likely owing to their ImageNet-21K pretraining and stronger inductive biases. The Swin Transformer is competitive, suggesting that global self-attention captures long-range contextual inconsistencies that

CNNs miss, though at the cost of higher memory footprint. MSA-CNN surpasses all baselines, despite not using an explicit frequency branch. We attribute this to the complementary role of multi-scale pyramid features in capturing spectral-like texture differences at different spatial frequencies: fine-resolution features (P1) encode high-frequency texture; coarser features (P2, P3) encode medium and low spatial frequencies. SE attention further sharpens these representations by focusing on channels most correlated with synthesis artefacts.

6.2 Confusion Matrix Analysis

Table 3 shows the averaged confusion matrix for MSA-CNN on the 436-volume test set. Of 217 truly manipulated volumes, 210 are correctly identified (7 false negatives). Of 219 truly authentic volumes, 213 are correctly classified (6 false positives). False negatives correspond primarily to nodule-removal attacks on small nodules (≥ 6 mm, 4 of 7 cases) and nodule injections at the lung periphery where boundary artefacts are less pronounced (3 of 7 cases). False positives are associated predominantly with volumes containing unusual parenchymal patterns (honeycombing, ground-glass opacities) that superficially resemble GAN texture statistics.

Table 3. Averaged confusion matrix for MSA-CNN on the test set (436 volumes, 5 seeds).

	Pred: Authentic	Pred: Manipulated
True: Authentic	TN = 213.0	FP = 6.0
True: Manipulated	FN = 7.0	TP = 210.0

6.3 Slice-Level vs. Volume-Level Performance

To understand the value of the majority voting aggregation step, Table 4 compares slice-level and volume-level AUC for all models. Volume-level aggregation consistently improves AUC by a mean of 1.8 percentage points across models. This is expected: individual slices far from the manipulated region contain little discriminative signal, and their inclusion at the slice level introduces noise. Majority voting suppresses this noise by averaging over the full volume.

Table 4. Slice-level vs. volume-level AUC-ROC on the test set (MSA-CNN, 5-seed mean).

Model	Slice AUC	Volume AUC	Δ AUC
ResNet-50	0.934	0.952	+0.018
EfficientNet-B5	0.954	0.971	+0.017
MSA-CNN (Ours)	0.973	0.991	+0.018

7. Ablation Study

We conduct a systematic ablation of MSA-CNN to quantify the contribution of each design element. Ablations are evaluated on the validation set (5-seed mean AUC). Table 5 summarises results.

Table 5. Component ablation study on the validation set (5-seed mean AUC and ACC).

Configuration	Val. AUC	Val. ACC (%)	Δ vs. Full Model
Full MSA-CNN	0.989	96.82	—
w/o Multi-Scale Pyramid (single-scale P3 only)	0.978	94.41	-0.011
w/o SE Attention (standard conv after FPN)	0.982	95.27	-0.007
w/o Auxiliary Heads	0.984	95.68	-0.005
w/o CutMix Augmentation	0.981	95.03	-0.008
w/o Random Erasing	0.986	96.14	-0.003
ResNet-50 \rightarrow EfficientNet-B4 Backbone	0.985	95.91	-0.004
FPN lateral width 128 (vs. 256)	0.987	96.38	-0.002
Dropout 0.3 (vs. 0.45)	0.983	95.44	-0.006

The most impactful component is the multi-scale pyramid (Δ AUC = -0.011 when removed), confirming that artefacts at multiple spatial frequencies are essential discriminative cues in this task. SE attention is the second most important element (Δ AUC = -0.007), consistent with the spatially sparse nature of manipulation artefacts. CutMix augmentation (Δ AUC = -0.008) is notably helpful, likely because mixing clean and manipulated regions during training forces the model to attend to local manipulation statistics rather

than holistic volume appearance. Auxiliary heads and random erasing provide smaller but consistent benefits. Replacing the ResNet-50 backbone with EfficientNet-B4 reduces validation AUC by 0.004, suggesting that the deeper residual feature hierarchy of ResNet-50 is marginally more suited to the FPN extraction setup we employ. Using a narrower FPN (128 channels) incurs only a 0.002 AUC penalty, suggesting that the architecture is not highly sensitive to FPN width beyond a minimum capacity. Dropout rate has a moderate effect; 0.45 outperforms 0.30 by 0.006 AUC, consistent with the moderate dataset size encouraging stronger regularisation.

8. Discussion

8.1 Clinical Implications

MSA-CNN’s sensitivity of 96.82% and specificity of 97.43% translate to meaningful clinical impact estimates. Assuming a manipulation prevalence of 0.1% (one manipulated study per 1,000), Bayes’ theorem yields a positive predictive value of approximately 3.1%—that is, only 3 in every 100 flagged studies would actually be manipulated. This low PPV in a prevalence-scarce environment highlights the necessity of not treating automated detection as a final arbiter. Instead, the practical use case is as a screening filter: studies flagged by MSA-CNN are routed to a secondary expert review queue, which dramatically reduces the fraction of manipulated scans that reach a radiologist undetected while keeping the manual review burden manageable. In a higher-prevalence adversarial context—such as a healthcare provider under active targeted attack, where prevalence might be estimated at 1%—PPV rises to approximately 26%, making automated detection substantially more actionable. We recommend that any clinical deployment of a system like MSA-CNN be accompanied by ongoing prevalence monitoring and threshold recalibration.

8.2 Limitations

Several limitations of this study must be acknowledged. First, the CT-GAN dataset derives exclusively from LIDC-IDRI, which was collected under a specific set of acquisition protocols. Generalisation to scanners, reconstruction kernels, and patient populations not represented in LIDC-IDRI has not been validated. Prospective evaluation on out-of-distribution data is an essential next step before clinical deployment.

Second, the label granularity in our experiments is at the volume level, following the released CT-GAN annotations. Fine-grained slice- or region-level localisation of the manipulated area—necessary for a radiologist to efficiently investigate a flagged study—is not addressed. Extending MSA-CNN to a detection-and-segmentation framework is a natural and important direction for future work.

Third, we evaluate only against the specific GAN architecture used in CT-GAN. More recent synthesis

approaches, including StyleGAN variants (Karras et al., 2020) and diffusion-based in-painting, may leave different statistical signatures. Our spectral-implicit multi-scale approach should generalise to some extent, but explicit cross-GAN generalisation experiments are warranted.

Finally, all architectures are tested under benign (non-adversarial) conditions. A motivated adversary aware of the detection system could craft adversarial perturbations specifically designed to fool the classifier. Investigating adversarial robustness and appropriate defences is a critical research direction that complements the detection accuracy results reported here.

8.3 Ethical Considerations

Research in medical deepfake detection carries an inherent dual-use tension: detailed documentation of the detection methodology could, in principle, inform adversaries seeking to evade detection. We have chosen to publish the detection method in full because the benefits of enabling other researchers to build on and improve upon our approach outweigh this risk, particularly given that the CT-GAN attack methodology itself is already public. We have not published our trained model weights without access controls; researchers requesting access for non-commercial purposes may contact the corresponding author. No patient data beyond what is included in the de-identified LIDC-IDRI collection was used or generated in this study.

9. Conclusion

This paper has presented a comprehensive deep learning investigation into the detection of GAN-synthesised lung CT images using the CT-GAN benchmark. We compared seven architectures under a rigorous common protocol and proposed MSA-CNN, which achieves 97.14% accuracy and an AUC-ROC of 0.991, outperforming all baselines with statistical significance. Ablation experiments confirm that multi-scale feature aggregation and squeeze-and-excitation attention are the most important architectural contributions, while CutMix augmentation provides the largest training-side benefit. Our results establish that automated medical deepfake detection in lung CT imaging is feasible at clinically relevant accuracy levels using current deep learning methodology. We hope that the benchmark protocol, architecture specifications, and dataset splits released alongside this paper will accelerate further progress in this urgent domain.

References

- [1] Armato, S. G., McLennan, G., Bidaut, L., McNitt-Gray, M. F., Meyer, C. R., Reeves, A. P., ... & Clarke, L. P. (2011). The Lung Image Database Consortium (LIDC) and Image Database Resource Initiative (IDRI): a completed reference database of lung nodules on CT scans. *Medical Physics*, 38(2), 915–931.
- [2] Rong C, Li Z, Li R, Wang Y. Spatial-aware contrastive learning for cross-domain medical image registration. *Med Phys*. 2024 Nov;51(11):8141-8150. doi: 10.1002/mp.17311. Epub 2024 Jul 19. PMID: 39031488.
- [3] Durall, R., Keuper, M., & Keuper, J. (2020). Watch your up-convolution: CNN based generative deep neural networks are failing to reproduce spectral distributions. *CVPR 2020*, 9930–9938.
- [4] Frank, J., Eisenhofer, T., Schönherr, L., Fischer, A., Kolossa, D., & Holz, T. (2020). Leveraging frequency analysis for deep fake image recognition. *ICML 2020*, 3270–3280.
- [5] Frid-Adar, M., Diamant, I., Klang, E., Amitai, M., Goldberger, J., & Greenspan, H. (2018). GAN-based synthetic medical image augmentation for increased CNN performance in liver lesion classification. *Neurocomputing*, 321, 321–331.
- [6] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., ... & Bengio, Y. (2014). Generative adversarial nets. *NeurIPS 2014*, 2672–2680.
- [7] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *CVPR 2016*, 770–778.
- [8] Hu, J., Shen, L., & Sun, G. (2018). Squeeze-and-excitation networks. *CVPR 2018*, 7132–7141.
- [9] Huang, G., Liu, Z., Van Der Maaten, L., & Weinberger, K. Q. (2017). Densely connected convolutional networks. *CVPR 2017*, 4700–4708.
- [10] Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., & Aila, T. (2020). Analyzing and improving the image quality of StyleGAN. *CVPR 2020*, 8107–8116.
- [11] Khalid, H., Tariq, S., Kim, M., & Woo, S. S. (2021). FakeBuster: a DeepFakes detection tool for video conferencing scenarios. *IUI 2021*, 75–77.
- [12] Ali A, Basha HA, Thanuja K, Puneet, Gupta SK, Kim S. Enhancing tumor deepfake detection in MRI scans using adversarial feature fusion ensembles. *Sci Rep*. 2025 Dec 9;16(1):1667. doi: 10.1038/s41598-025-31231-7. PMID: 41366279; PMCID: PMC12800179.
- [13] Li, L., Bao, J., Zhang, T., Yang, H., Chen, D., Wen, F., & Guo, B. (2020). Face X-ray for more general face forgery detection. *CVPR 2020*, 5001–5010.
- [14] Lin, T. Y., Dollár, P., Girshick, R., He, K., Hariharan, B., & Belongie, S. (2017). Feature pyramid networks for object detection. *CVPR 2017*, 2117–2125.
- [15] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., ... & Guo, B. (2021). Swin Transformer: hierarchical vision transformer using shifted windows. *ICCV 2021*, 10012–10022.

- [16] Liu, Z., Mao, H., Wu, C. Y., Feichtenhofer, C., Darrell, T., & Xie, S. (2022). A ConvNet for the 2020s. *CVPR 2022*, 11976–11986.
- [17] Loshchilov, I., & Hutter, F. (2019). Decoupled weight decay regularization. *ICLR 2019*.
- [18] Mirsky, Y., Mahler, T., Shelef, I., & Elovici, Y. (2019). CT-GAN: malicious tampering of 3D medical imagery using deep learning. *USENIX Security Symposium 2019*, 461–478.
- [19] Motamed, S., Rogalla, P., & Khalvati, F. (2021). Data augmentation using Generative Adversarial Networks (GANs) for GAN-based detection of Pneumonia and COVID-19 in chest X-ray images. *Informatics in Medicine Unlocked*, 27, 100779.
- [20] Nataraj, L., Mohammed, T. M., Manjunath, B. S., Chandrasekaran, S., Flenner, A., Bappy, J. H., & Roy-Chowdhury, A. K. (2019). Detecting GAN generated fake images using co-occurrence matrices. *Electronic Imaging 2019*, 532-1–532-7.
- [21] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. *CVPR 2022*, 10684–10695.
- [22] Rossler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., & Niessner, M. (2019). FaceForensics++: learning to detect manipulated facial images. *ICCV 2019*, 1–11.
- [23] Salehinejad, H., Valaee, S., Khalvati, F., Merat, S., & Salehinejad, E. (2018). Generalization of deep neural networks for chest pathology classification in X-rays using generative adversarial networks. *ICASSP 2018*, 990–994.
- [24] Yun, S., Han, D., Oh, S. J., Chun, S., Choe, J., & Yoo, Y. (2019). CutMix: training strategy that makes use of sample mixing and its data augmentation. *ICCV 2019*, 8024–8033.
- [25] Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2021). Score-based generative modeling through stochastic differential equations. *ICLR 2021*.
- [26] Tan, M., & Le, Q. V. (2019). EfficientNet: rethinking model scaling for convolutional neural networks. *ICML 2019*, 6105–6114.
- [27] Wolterink, J. M., Dinkla, A. M., Savenije, M. H. F., Seevinck, P. R., van den Berg, C. A. T., & Isgum, I. (2017). Deep MR to CT synthesis using unpaired data. *MICCAI 2017*, 14–23.