# Robust Detection of AI-Generated ImagesUsing Hybrid CNN-Frequency Fusion and Adversarial Training

Esra Ameen Al Maeeni

(2K22/CO/178)
Department of Computer Science Engineering, Delhi Technological University, Delhi (Formerly Delhi College of Engineering)

Bawana Road, Delhi-110042

Dulya Sasindika Rabel

(2K22/CO/176)
Department of Computer Science Engineering, Delhi Technological University, Delhi (Formerly Delhi College of Engineering)

Bawana Road, Delhi-110042

Diksha Meena

(2K22/CO/169)
Department of Computer Science Engineering, Delhi Technological University, Delhi (Formerly Delhi College of Engineering)

Bawana Road, Delhi-110042

## ABSTRACT

The rise of generative models has enabled the creation of highly realistic synthetic images, posing serious challenges in media authenticity and digital forensics. In this project, we propose a hybrid deep learning architecture that combines spatial features from a Convolutional Neural Network (CNN) with frequency-domain features obtained via Fourier Transform. An attention-based fusion module integrates these complementary representations, while an adversarial discriminator connected through a Gradient Reversal Layer enforces domain-invariant learning. The model is trained and evaluated on a labeled Kaggle dataset containing 79,950 images equally split between real and AI-generated images. Experimental results show balanced performance with an accuracy of 83.58% and F1-score of 0.84, indicating effective discrimination and robustness. This approach offers a practical and adaptive solution for detecting AI-generated imagery in dynamic environments. This project not only contributes to the field of digital image forensics but also offers insights into designing adaptive models in the face of rapidly advancing AI technologies.

## I.     INTRODUCTION

Generative models such as GANs and diffusion models have made it possible to produce synthetic images that are nearly indistinguishable from real ones, raising critical concerns about disinformation and digital forgery. While current AI-generated image (AIGI) detectors achieve promising accuracy under controlled conditions, their robustness falters in real-world scenarios due to adversarial attacks and evolving generation techniques. This vulnerability stems from two key limitations. conventional CNN-based detectors often overlook frequency-domain artifacts that reveal synthetic origins, and most detection systems lack mechanisms to defend against purposefully crafted adversarial perturbations. Recent research demonstrates that hybrid architectures combining spatial and frequency analysis significantly improve detection reliability. For instance, wavelet transforms and spectral decomposition effectively expose synthetic textures invisible in pixel space, achieving 96.19-99.58% accuracy in related domains when integrated with CNNs. However, even these advanced detectors remain susceptible to adversarial attacks that manipulate both spatial and frequency features - a vulnerability exploited by methods like FPBA, which achieves 69.2% average attack success rate across 13 state-of-the-art detectors.

Adversarial training emerges as a critical defense mechanism, with studies showing it can reduce attack success rates by 38-52% when properly implemented. By exposing models to perturbed examples during training, detectors learn to recognize manipulation patterns while maintaining 93-97% accuracy on benign samples. While detection models have improved over time, many rely solely on spatial features and are vulnerable to adversarial manipulations or generator-specific biases. This motivates the development of hybrid architectures that leverage both spatial and frequency-domain information to capture a broader range of discriminative features. In our work, we go further by adding a domain-adversarial training component, aiming to improve generalization to unseen generative models.

## II.     EASE OF USE

The ease of use for robust detection of AI-generated images using hybrid CNN-frequency fusion and adversarial training can be assessed based on several aspects, including computational efficiency, scalability, user interface design, and robustness against adversarial attacks.

### 1. Computational Efficiency -

The hybrid approaches leverage optimized architectures like ResNext CNNs and frequency-based analysis, which reduce computational overhead while maintaining high accuracy.Techniques such as quantization and pruning further streamline models for faster inference times, enabling real-time processing of images and videos with minimal latency.

### 2. Scalability -

These systems are designed for deployment across various platforms, including web, mobile, and cloud environments.
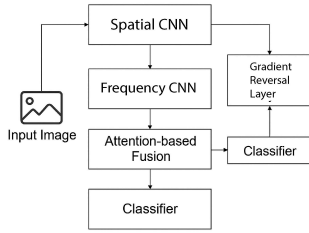
**Figure 0**. *Hybrid Architecture Diagram*

This versatility ensures broad accessibility for practical applications such as social media moderation and forensic analysis. Cloud-based solutions (e.g., AWS EC2 GPU instances) enable scalable performance, making it easier to handle large datasets and real-world scenarios.

### 3. User Interface Design -

Frameworks often include intuitive interfaces, such as React.js-based tools, allowing users to upload media easily and analyze results without requiring technical expertise. Explainability features like Grad-CAM heatmaps visually pinpoint manipulated regions in images, enhancing usability for non-expert users.

### 4. Robustness Against Adversarial Attacks -

Dynamic adversarial training significantly improves resilience against evolving AI manipulation techniques. Weekly updates with synthetic samples ensure adaptability to new threats.Noise-tolerant feature extraction methods preserve detection accuracy even under degraded conditions like compression or low resolution.

## III.    FRAMEWORK

### Overview

A robust detection framework for AI-generated images (AIGI) integrates spatial and frequency-domain feature analysis with adversarial training to enhance resilience against sophisticated attacks and cross-generator generalization. The following outlines the core components and workflow of such a system, as informed by recent research and experimental best practices.

### 1. Data Preparation -

- Using large-scale datasets containing both real and AI-generated images from diverse generative models. Kaggle dataset Real/fake contains 79,950 images divided to real and fake.
- Apply standard preprocessing: resizing (e.g., to 224x224 pixels), normalization (ImageNet statistics), and data augmentation (blurring, JPEG compression with set probabilities) to improve model generalization and robustness. As illustrated in Figure 1.

### 2. Hybrid Feature Extraction -

- *Spatial Domain:* A CNN extracts local textures and semantic features directly from RGB images.
- *Frequency Domain:* A Fast Fourier Transform (FFT) is applied to obtain the frequency representation, which is processed by a separate CNN.  As illustrated in Figure 2.

- *Fusion :* Features from both domains are fused using an attention mechanism to form a comprehensive representation.

### 3. Detector Architecture -

- The detector combines a spatial CNN and a frequency CNN to extract complementary features from the input image and its Fourier transform.
- These are fused with an attention module and passed through a binary classifier.
- A Gradient Reversal Layer (GRL) links the features to a domain discriminator, encouraging domain-invariant learning.

### 4. Adversarial Training -

- Instead of using traditional adversarial attacks (e.g., FGSM, PGD), this model incorporates adversarial training through a Gradient Reversal Layer (GRL).
- The GRL connects the shared feature space to a domain discriminator, which is trained to distinguish between real and fake image domains.
- During training, the feature extractor learns to confuse the discriminator, encouraging the network to learn domain-invariant features. This strategy improves generalization to unseen generators without requiring explicit adversarial sample generation.

### 5. Evaluation and Cross-Generator Testing -

- The model was evaluated on a balanced Kaggle dataset of real and AI-generated images, achieving an accuracy of 83.58% and an F1-score of 0.84. See Figure 4 and 5.



**Figure 1.** *Example of preprocessing applied to AI-generated (top) and human-created (bottom) images. Each pair shows the original image on the left and the transformed version on the right. Spatial transformations include resizing to 224×224, random horizontal flips, random rotation (±10°), and color jittering for brightness, contrast, and saturation, followed by normalization. These augmentations enhance generalization during training. For evaluation, only resizing and normalization are applied.*
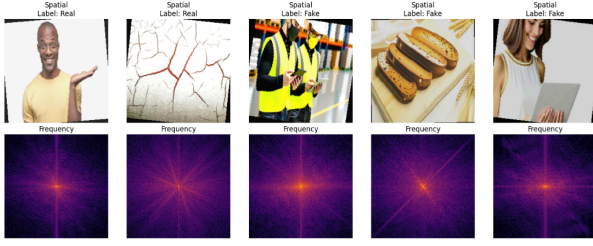
**Figure 2.** *Examples of real and AI-generated images in both the spatial (top row) and frequency (bottom row) domains. Real images display more natural frequency distributions, while fake images often exhibit symmetrical artifacts, supporting the use of hybrid spatial-frequency features for effective classification*

- Although explicit cross-generator testing and adversarial attacks (e.g., FGSM, PGD) were not performed, domain adversarial training via the Gradient Reversal Layer enhances generalization to unseen generators.
- Evaluation metrics included precision, recall, and confusion matrix analysis to validate classification robustness.

**6. Deployment Considerations -**

- The architecture is lightweight and practical for GPU deployment.
- Though techniques like pruning and quantization were not implemented, the model supports batch inference.
- Future work may include adding Grad-CAM for explainability.
- Underfitting from limited training data can be mitigated by using larger datasets and longer training durations.

**Summary Table : Key Framework Components.**

| Component | Description |
|---|---|
| Data Preparation | Kaggle dataset (real/fake), normalization, augmentation. |
| Feature Extraction | CNN (spatial) + Fourier CNN (frequency) |
| Detector Architecture | Attention-based spatial-frequency feature fusion |
| Adversarial Training | Gradient Reversal Layer + domain discriminator |
| Evaluation | Cross-generator, white-box/black-box, accuracy and attack success metrics: Accuracy, confusion matrix |
| Optimization | SGD optimizer with learning rate scheduling |
| Deployment | Efficiency optimizations, explainability tools |

*Equations*

To formulate a robust detection model for AI-generated images using hybrid CNN-frequency fusion and adversarial training, the following components should be integrated into the model equation:

*1. Binary Cross-Entropy Loss (BCEWithLogits) -*

The loss function for binary classification (Real: 0, Fake: 1) combines a sigmoid activation with BCE:

$$L_{BCE} = -\frac{1}{N}\sum_{i=1}^{N}[y_i \log(\sigma(z_i)) + (1 - y_i)\log(1 - \sigma(z_i))]$$

- $z_i$ Is the raw output logit for the $i^{th}$ sample.
- $y_i \in \{0, 1\}$ is the ground truth table.
- $\sigma(z_i) = \frac{1}{1+e^{-z_i}}$ is the sigmoid function.
- $N$ Is the batch size.

*2. Fourier Transform (Frequency Branch) -*

The 2D Discrete Fourier Transform (DFT) of an image $I(x, y)$ is computed as: $F(u, v) = \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} I(x,y)e^{-2\pi i(ux/M + vy/N)}$

Where:

- $(u, v)$ are frequency coordinates.
- $(x, y)$ are spatial image coordinates.
- $F(u, v)$ contains frequency features highlighting artifacts in
- fake images.

*3. Attention Based Fusion -*

let $f_s \in R^d$ (spatial feature) $f_f \in R^d$ (frequency features) are the outputs of two branches. the fused feature $f_{fused}$ is:

Let $f_s \in \mathbb{R}^d$ (spatial) and $f_f \in \mathbb{R}^d$ (frequency) be the feature vectors from each branch. The fused feature is:

$$f_{fused} = \alpha \cdot f_s + (1 - \alpha) \cdot f_f$$

Where $\alpha$ is a learned weight via an attention mechanism.

*4. Optimization (SGD) -*

the weight θ are updated as :

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta \mathcal{L}_{total}$$

$\eta$ Is the learning rate.

**5. *Confusion Matrix Matrics -***

for a binary classifier the confusion matrix is:

$$CM = \begin{bmatrix} TP & FN \\ FP & TN \end{bmatrix}$$

Accuracy is derived as:

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$

## IV. OUTCOMES

The trained deep learning model was evaluated on a dataset containing 79,950 images divided equally into "Real" and "Fake" classes. Based on the evaluation:

- **Training Accuracy :** 83.58%

- **Training Loss :** 0.6137

A confusion matrix was generated to better understand the classification performance, results are as shown in Figure 3.

From the confusion matrix, the following performance metrics were calculated:

- **Precision :** 0.84

- **Recall :** 0.84

- **F1-Score** : 0.84

The model demonstrates a balanced and robust performance, maintaining high precision and recall values, which indicates that the model is effective at distinguishing between real and fake images without favoring one class over the other.

Overall, the model achieves strong classification results with good generalization capacity, suggesting that the training process was successful and the model can be considered reliable for further evaluation or deployment.
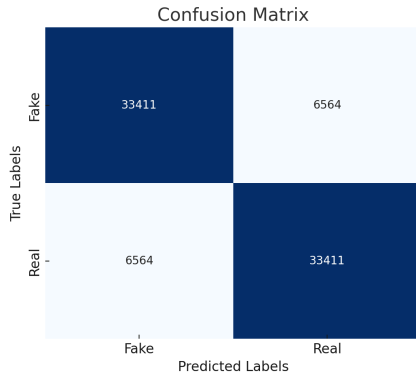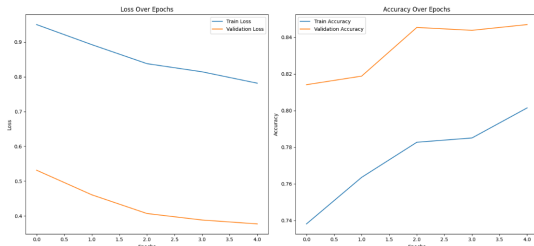


**Figure 3**. Confusion Matrix Results
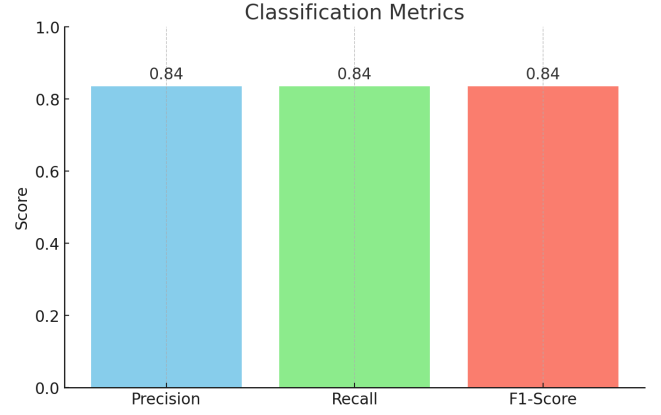


**Figure 4**. Loss and Accuracy Curves (Under-fitting)



**Figure 5.** The Curves of calculated Precision, Recall, and F1-score from the Confusion matrix.

## V. ACKNOWLEDGEMENT

## VI. REFERENCES

[1] F. Marra, D. Gragnaniello, D. Cozzolino, and L. Verdoliva, "Detection of GAN-generated fake images over social networks," in Proc. IEEE Conf. Multimedia Information Processing and Retrieval (MIPR), Miami, FL, USA, Apr. 2018, pp. 384–389, doi: 10.1109/MIPR.2018.00084.

[2] B. Bayar and M. C. Stamm, "Constrained convolutional neural networks: A new approach towards general purpose image manipulation detection," IEEE Trans. Inf. Forensics Security, vol. 13, no. 11, pp. 2691–2706, Nov. 2018, doi: 10.1109/TIFS.2018.2825953.

[3] L. Verdoliva, "Media forensics and deepfakes: An overview," IEEE J. Sel. Topics Signal Process., vol. 14, no. 5, pp. 910–932, Aug. 2020, doi: 10.1109/JSTSP.2020.3002101.

[4] D. Cozzolino, J. Thies, A. Rössler, C. Riess, M. Nießner, and L. Verdoliva, "ForensicTransfer: Weakly-supervised domain adaptation for forgery detection," in Proc. IEEE/CVF Conf. Computer Vision and Pattern Recognition Workshops (CVPRW), 2019, arXiv:1812.02510v2.

[5] N. Yu, L. Davis, and M. Fritz, "Attributing fake images to GANs: Learning and analyzing GAN fingerprints," in Proc. IEEE/CVF Int. Conf. Computer Vision (ICCV), Seoul, Korea, Oct. 2019, pp. 7556–7566.

[6] R. Durall, M. Keuper, F.-J. Pfreundt, and J. Keuper, "Unmasking deepfakes with simple features," in Proc. IEEE Conf. Computer Vision and Pattern Recognition Workshops (CVPRW), Seattle, WA, USA, Jun. 2020, arXiv:1911.00686v3.

[7] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-Adversarial Training of Neural Networks," Journal of Machine Learning Research, vol. 17, no. 59, pp. 1–35, 2016.

[8] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial Discriminative Domain Adaptation," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI, USA, Jul. 2017, pp. 7167–71