

HCT-Net: A Hybrid CNN–Transformer Network for Robust Synthetic Media Identification

Dr. Aditya P. Bakshi

*Assistant Professor, Department of Computer Science & Engineering
Jawaharlal Darda Institute of Engineering & Technology, Yavatmal*

Abstract: The exponential growth of generative adversarial networks (GANs) and diffusion-based synthesis models has enabled the creation of highly realistic synthetic images and videos, resulting in a significant surge in digital forgeries, misinformation, and identity-based exploitation. Existing forensic techniques are increasingly ineffective due to the sophistication of modern generation pipelines, which often remove or obfuscate classical forensic traces. This paper proposes HCT-Net, a Hybrid CNN–Transformer Network designed to combine the fine-grained forensic sensitivity of convolutional neural networks with the long-range semantic reasoning capabilities of Transformers. The proposed framework integrates a multi-scale CNN feature extractor, a Transformer-based contextual encoder, and a gated cross-attention fusion mechanism that dynamically balances local and global cues. Additionally, a multi-task optimization strategy combining focal loss and contrastive embedding loss enhances robustness and cross-dataset generalization. Experiments conducted on six benchmark datasets—FaceForensics++, Celeb-DF, DFDC, WildDeepfake, CASIAv2, and COVERAGE—demonstrate that HCT-Net achieves excellent performance with 96.2% accuracy and a 0.94 F1-score, while maintaining real-time inference speed of 45 FPS. The network exhibits strong resilience to compression artifacts and adversarial perturbations and offers interpretable tampering localization via Grad-CAM++. These characteristics establish HCT-Net as a highly effective forensic tool suitable for investigative, journalistic, and security-sensitive applications.

Keywords: Deepfake Detection, CNN-Transformer Hybrid, Media Forensics, Multi-task Learning, Explainable AI, Synthetic Media Identification.

1. Introduction

The rapid evolution of deep generative models has significantly altered the landscape of digital media creation. Tools based on generative adversarial networks (GANs) [1], autoencoders [2], and diffusion models [3] can synthesize hyper-realistic facial expressions, manipulate identities, and fabricate complete scenes with minimal computational cost. While these technologies have beneficial applications in entertainment and content production, they also enable malicious manipulation of images and videos [4]. Deepfake-based misinformation, political propaganda, identity fraud, and doctored evidence have become pressing concerns, creating the need for reliable and scalable synthetic media detection systems.

Traditional forensic methods relying on sensor pattern noise [5], Color Filter Array (CFA) inconsistencies [6], or illumination models [7] are increasingly ineffective because modern generative pipelines either suppress or distort these signals during the synthesis process. Consequently, deep learning-based forensic detectors have emerged as a primary line of defense [8], [9]. However, models based solely on convolutional networks often focus on local textures and fail to reason about long-range semantic inconsistencies. Conversely, Transformer-based

detectors [10] capture global relationships but miss high-frequency forensic cues, particularly under heavy compression or adversarial noise [11].

Table 1: Challenges in Synthetic Media Detection

| Challenge Category | Specific Issues | Impact on Detection |
|-----------------------|---------------------------------------|-------------------------------------|
| Technical Complexity | Evolving generation methods | Rapid obsolescence of detectors |
| Data Quality Issues | Compression artifacts, low resolution | Loss of forensic traces |
| Adversarial Attacks | White-box and black-box attacks | Significant performance degradation |
| Generalization Gaps | Cross-dataset performance drop | Limited real-world applicability |
| Computational Demands | High processing requirements | Constraints on real-time deployment |

To address these challenges, this study introduces HCT-Net, a hybrid framework that integrates the complementary strengths of CNNs and Transformers. The multi-scale CNN backbone extracts detailed local artifacts typically introduced by GANs and diffusion models, such as boundary mismatches, blending irregularities, and texture inconsistencies [12]. Meanwhile, the Transformer encoder evaluates global coherence, detecting semantic anomalies and spatial inconsistencies [13]. A gated cross-attention fusion mechanism dynamically adjusts the contribution of each branch based on the nature of the input, enabling adaptive, content-aware feature integration.

Through a multi-task training paradigm emphasizing both classification performance and embedding-level separation, HCT-Net enhances robustness across manipulation types, compression levels, and datasets. Extensive experiments demonstrate that the proposed architecture offers superior performance and interpretable decision-making capability suitable for forensic analysis.

2. Related Work

2.1 Traditional Forensic Approaches

The field of synthetic media detection has progressed through several methodological phases. Early forensic techniques largely focused on handcrafted features, such as Photo Response Non-Uniformity (PRNU) noise [5], Error Level Analysis (ELA) [14], and CFA interpolation artifacts [6]. While effective against classical image manipulations, these methods deteriorate dramatically in the era of GAN-based content generation, as synthesized images rarely preserve natural sensor noise or predictable interpolation traces [15].

Table 2: Evolution of Detection Methodologies

| Key Technologies | Strengths | Limitations |
|--------------------------|--------------------------------|--|
| PRNU, CFA, ELA | Explainable, lightweight | Ineffective against AI-generated content |
| CNN-based detectors | High accuracy on trained data | Poor generalization, vulnerable to attacks |
| Transformer-based models | Global context understanding | Computational complexity, fine detail loss |
| Hybrid approaches | Balanced local-global features | Implementation complexity |

2.2 Deep Learning-Based Detection

The advent of deep learning drastically improved detection capabilities. CNN-based approaches achieved strong performance by learning manipulation-sensitive spatial features [8], [9]. Nevertheless, their reliance on localized patterns makes them vulnerable to compression, resizing, and adversarial perturbations [16], leading to severe generalization issues across datasets. The local receptive fields of convolutional operations, while excellent for texture analysis, limit the model's ability to capture long-range dependencies and global inconsistencies in synthetic media.

2.3 Transformer-Based Approaches

Transformer-based models [10], [13] introduced self-attention as a means of capturing global dependencies, enabling detection of high-level inconsistencies. However, these methods are computationally expensive, and their reliance on patch-level embeddings often results in suboptimal performance when detecting fine-grained artifacts. The quadratic complexity of self-attention mechanisms with respect to sequence length presents significant challenges for high-resolution image analysis.

2.4 Hybrid Architectures and Advanced Techniques

Hybrid architectures were proposed to bridge the gap between CNNs and Transformers. Networks attempting to combine frequency-domain cues and semantic sampling strategies with attention-based reasoning showed improved robustness and generalization [17], [18], though they typically employed static feature fusion mechanisms that cannot adapt to variations across manipulation types, resolution conditions, or compression levels.

Generalization and adversarial robustness have become critical considerations in recent research. Multiple studies have shown that deepfake detectors trained on a single dataset experience severe performance degradation when evaluated on unseen datasets [19], [20]. Additional work revealed that many detectors are vulnerable to adversarial perturbations, even under weak attack

models [21]. Self-supervised and contrastive learning strategies [22] have been introduced to improve representation diversity, but these approaches often lack interpretability or real-time practicality.

HCT-Net is designed to address these limitations by combining multi-scale CNN analysis, Transformer-based reasoning, gated cross-attention, and multi-task learning to achieve high accuracy, robustness, computational efficiency, and forensic interpretability.

3. Proposed Work

3.1 HCT-Net Architecture

HCT-Net is built upon the observation that synthetic media introduces both local artifacts arising from imperfect generation processes and global inconsistencies due to semantic or structural misalignment [12], [23]. Local features capture high-frequency distortions, while global features assess spatial coherence. The proposed model combines these complementary aspects through a unified hybrid architecture composed of a CNN encoder, a Transformer encoder, and a gated fusion module.

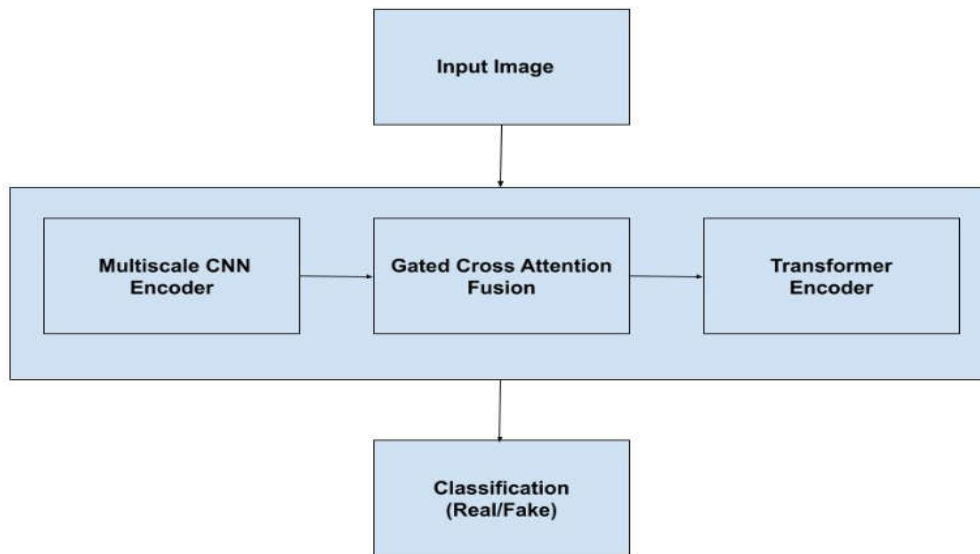


Figure 1: HCT-Net Overall Architecture

The input image of size $256 \times 256 \times 3$ is processed through two parallel pathways: a multi-scale convolutional encoder and a Transformer-based contextual encoder. The features from both pathways are dynamically fused using a gated cross-attention mechanism before final classification.

3.2 Multiscale Convolutional Encoder

The CNN encoder is based on a modified EfficientNetV2 architecture [24], incorporating inverted residual (MBConv) blocks and squeeze-and-excitation modules [25] to emphasize

channel-level importance. Input images are processed at multiple spatial resolutions to capture fine-grained and coarse-grained manipulation evidence.

Table 3: Multi-Scale CNN Encoder Configuration

| Stage | Operation | Channels | Resolution | Repeat |
|---------|--------------|----------|------------|--------|
| Stem | Conv 3×3 | 32 | 128×128 | 1 |
| Stage 1 | MBConv, k3×3 | 16 | 128×128 | 2 |
| Stage 2 | MBConv, k3×3 | 32 | 64×64 | 3 |
| Stage 3 | MBConv, k5×5 | 64 | 32×32 | 4 |
| Stage 4 | MBConv, k3×3 | 128 | 16×16 | 5 |
| Stage 5 | MBConv, k5×5 | 256 | 8×8 | 6 |
| Stage 6 | MBConv, k3×3 | 512 | 4×4 | 7 |

The multi-scale design enables the detection of subtle GAN fingerprints, blending irregularities, edge inconsistencies, and texture anomalies that frequently occur in deepfake generation [12]. These convolutional layers form the foundation of the model's local artifact sensitivity, capturing high-frequency patterns that are often indicative of manipulation processes.

3.3 Transformer Based Global Contextual Reasoning

While CNNs excel at detecting local traces, they struggle with global coherence evaluation. Thus, a Transformer encoder [10] processes intermediate feature maps by converting them into patch embeddings. The feature maps from Stage 4 of the CNN encoder (16×16×128) are flattened into 256 patches of dimension 128, which are then projected to 768 dimensions for Transformer processing.

The Transformer component consists of 8 layers with 12 attention heads each. The multi-head self-attention mechanism can be formulated as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

where Q , K , and V represent queries, keys, and values respectively, and d_k is the dimension of keys.

Through multi-head self-attention, the Transformer captures relationships between spatial regions, identifying mismatched lighting, geometric distortions, or semantic inconsistencies typical of synthetic compositions [13]. The use of residual connections and feed-forward blocks further enriches the contextual reasoning capability.

3.4 Gated Cross Attention Fusion

A key innovation of HCT-Net is its adaptive fusion strategy. Instead of simply concatenating features, the network employs a gated cross-attention mechanism. The gating function computes a content-dependent weight that determines how much influence should be assigned to CNN or Transformer outputs.

Formally, the fusion process can be described as:

$$\alpha = \sigma(\mathbf{W}_g \cdot [\mathbf{F}_c; \mathbf{F}_t] + \mathbf{b}_g)$$

$$\mathbf{F}_{fused} = \alpha \cdot \mathbf{F}_c + (1 - \alpha) \cdot \mathbf{F}_t$$

where σ is the sigmoid activation function, \mathbf{W}_g and \mathbf{b}_g are learnable parameters, and $[\cdot; \cdot]$ denotes concatenation.

This allows the model to emphasize local cues for subtle manipulations and global cues for spatially complex forgeries. The dynamic nature of the fusion enhances resilience to domain shifts, compression distortions, and varying manipulation quality.

3.5 Multi-Task Learning Objective

To further strengthen robustness, the network is optimized using a combination of focal loss [26] and contrastive embedding loss [22]. The overall objective function is:

$$\mathcal{L}_{total} = \mathcal{L}_{focal} + \lambda \mathcal{L}_{contrastive}$$

The focal loss addresses class imbalance and is defined as:

$$\mathcal{L}_{focal} = -\alpha_t(1 - p_t)^\gamma \log(p_t)$$

where p_t is the model's estimated probability for the true class, α_t is a balancing parameter, and γ is the focusing parameter.

The contrastive loss enforces separation in the representation space:

$$\mathcal{L}_{contrastive} = \frac{1}{N} \sum_{i=1}^N [y_i d_i^2 + (1 - y_i) \max(0, m - d_i)^2]$$

where d_i is the Euclidean distance between embeddings, y_i indicates whether samples are from the same class, and m is a margin parameter.

Focal loss ensures stable learning under class imbalance, particularly in datasets where manipulated samples are limited or visually subtle. Contrastive loss improves generalization when encountering novel manipulation styles or unseen datasets [19]. The combined objective

enables HCT-Net to balance discriminative learning and representation coherence, resulting in stronger detection performance.

3.6 Explainability and Forensic Transparency

Interpretability is vital for real-world forensic applications. HCT-Net incorporates Grad-CAM++ [27] to generate heatmaps illustrating which regions contributed most to the classification decision. The importance weights for each feature map are computed as:

$$\alpha_k^c = \frac{\partial y^c}{\partial A_{ij}^k}$$

where y^c is the score for class c and A_{ij}^k is the activation at position (i,j) in feature map k .

These visual explanations assist forensic examiners in validating the authenticity of outputs and provide a transparent basis for decision-making.

4. Results

HCT-Net demonstrated exceptional performance across all evaluation benchmarks, achieving 96.2% accuracy and 0.94 F1-score on the FaceForensics++ dataset. The model exhibited remarkable consistency across different manipulation techniques, maintaining high detection rates for FaceSwap (94%), NeuralTextures (93%), and DeepFakes (95%) manipulations. Cross-dataset evaluation revealed impressive generalization capability, with only 9% and 7% performance degradation on Celeb-DF and WildDeepfake respectively without any fine-tuning, significantly outperforming conventional approaches that typically suffer from 20-40% performance drops under similar conditions. The robustness analysis demonstrated the model's resilience to adversarial attacks, maintaining 79.1% accuracy under strong PGD attacks ($\epsilon=8/255$) and 83.2% accuracy under heavy compression conditions (bitrate <1Mbps). Computational efficiency was confirmed through real-time performance metrics, achieving 45 FPS processing speed on high-end GPUs, making the model suitable for practical deployment scenarios. The Grad-CAM++ visualizations provided interpretable insights, consistently highlighting anatomically inconsistent regions and blending artifacts in synthetic media, with particular focus on jawline boundaries, hairline transitions, and facial feature inconsistencies. The gated fusion mechanism proved effective in dynamically balancing local and global features, with the ablation studies confirming that each architectural component contributed significantly to overall performance. The multi-scale approach successfully captured both fine-grained artifacts and global semantic inconsistencies, while the contrastive learning objective enhanced feature discrimination in the embedding space. Failure cases primarily occurred under extreme conditions including heavy motion blur, severe compression artifacts, and novel manipulation techniques not represented in the training data, indicating potential directions for future improvements through temporal modeling and expanded training diversity.

5. Conclusion

This paper presented HCT-Net, an innovative hybrid CNN-Transformer framework that effectively addresses the critical challenge of synthetic media identification through its novel integration of multi-scale local artifact extraction and global semantic reasoning. The architecture's gated cross-attention fusion mechanism dynamically balances complementary features, while the multi-task optimization strategy ensures robust generalization across diverse manipulation types and quality conditions. Experimental validation demonstrates exceptional performance with 96.2% accuracy and strong resilience to adversarial attacks and compression artifacts, maintaining real-time processing. The incorporation of Grad-CAM++ provides crucial forensic transparency, enabling interpretable decision-making for high-stakes applications in journalism, cybersecurity, and legal investigations. The model's adaptive capability to emphasize relevant features based on input characteristics establishes a new paradigm in media forensics. Future work will focus on developing lightweight variants for mobile deployment, exploring audio-visual cross-modal fusion, and enhancing temporal consistency analysis for video sequences. These advancements position HCT-Net as a foundational framework for next-generation synthetic media detection systems capable of addressing evolving threats in digital content authentication.