# DeepTrace: Visual Signature Graphs for AI-Generated Image Detection Using Multi-Dimensional Feature Analysis and Graph-Based Anomaly Detection

Anagha Pradeep*
*Indiana University, Indianapolis, USA
Email: anaprad@iu.edu

*Abstract*—The proliferation of sophisticated generative models has created unprecedented challenges for image authenticity verification. This paper presents DeepTrace, a comprehensive, model-agnostic framework for detecting AI-generated images through multi-dimensional feature extraction combined with graph-based relationship modeling and machine learning classification. Our approach extracts eight complementary visual signatures capturing edge characteristics, frequency domain properties, noise statistics, and texture patterns. These features are mapped to an 8-dimensional space where similarity relationships are modeled via k-nearest neighbor graphs. We employ density-based clustering (DBSCAN) and isolation-based anomaly detection to reveal population-level patterns. Critically, our ultra-aggressive memory optimization enables processing of 100,000+ images with peak memory usage under 1.2 GB—a 33× reduction compared to naive approaches. We train a Random Forest classifier on the top 5 most discriminative features, achieving 85.65% accuracy, 84.48% precision, and 87.34% recall on 20,000 held-out test images. Extensive evaluation on balanced datasets (50,000 AI-generated from Stable Diffusion/DALL-E/Midjourney, 50,000 real images) demonstrates robust generalization across multiple generative models. Feature importance analysis reveals that FFT High-Frequency Ratio (34.2%) and Residual Kurtosis (31.1%) are most discriminative. Our unsupervised feature extraction combined with supervised classification provides interpretable, scalable, and practically deployable detection with clear forensic meaning for each decision signal.

*Index Terms*—AI-generated image detection, visual signatures, graph-based analysis, machine learning, feature extraction, generative models, image forensics, random forest classification

## I. INTRODUCTION

The emergence of diffusion models (Stable Diffusion, DALL-E), advanced GANs, and transformer-based image generators has fundamentally altered digital media authenticity. These models produce photorealistic images virtually indistinguishable from authentic photographs by human observers and pose severe risks including misinformation campaigns, deepfakes, copyright violations, and financial fraud.

Existing detection approaches suffer from critical limitations. Binary classifiers trained on specific generative models exhibit catastrophic performance degradation when encountering unseen models, with accuracy dropping from 98% to 53%. Traditional forensic methods exploit camera sensor fingerprints and compression artifacts irrelevant for synthetic images. Deep learning approaches require extensive labeled datasets, lack interpretability, and are computationally expensive.

DeepTrace addresses these gaps through:

1) **Model-Agnostic Detection**: Works across Stable Diffusion, DALL-E, Midjourney without retraining
2) **Comprehensive Features**: Eight complementary visual signatures capturing distinct generative artifacts
3) **Scalable Architecture**: Processes 100K+ images with memory optimization achieving 33× reduction
4) **High Accuracy**: Random Forest classification achieves 85.65% accuracy on balanced test set
5) **Interpretable Decisions**: Each feature has clear forensic meaning; feature importance quantified
6) **Production Ready**: Uses only standard libraries, runs on consumer hardware

## II. RELATED WORK

### A. Image Forensics

Traditional methods exploit PRNU (Photo Response Non-Uniformity) as camera fingerprints and JPEG compression artifacts. These techniques work well for in-distribution images but fail completely on synthetic content designed to avoid detection. Farid's foundational work established key artifacts but assumes authentic camera imaging pipelines.

### B. Generative Model Artifacts

Wang et al. demonstrated that CNNs achieve 99.8% accuracy detecting GAN artifacts on in-distribution data but degrade to 53% on out-of-distribution models. GANs exhibit spectral anomalies, boundary artifacts, mode collapse, and color shifts. Diffusion models create smooth transitions, high-frequency noise patterns, and characteristic texture uniformity from iterative denoising fundamentally different artifact profiles than GANs.

## C. Deep Learning Detection

CNN-based classifiers fine-tune ResNet/EfficientNet architectures, achieving 95-98% accuracy on balanced test sets. Vision Transformers show modest improvement (80-85% cross-model) but still require model-specific fine-tuning. Ensemble approaches provide robustness but increase computational cost and deployment complexity.

## D. Graph-Based and Unsupervised Methods

Graph Neural Networks enable anomaly detection through spectral methods and message-passing. Classical graph analysis combined with comprehensive feature extraction enables interpretable anomaly detection without supervision. Zhou et al.'s comprehensive review demonstrates effectiveness for node-level anomaly detection through structural context.

## III. METHOD

### A. System Architecture

DeepTrace operates through five sequential stages:
**Stage 1: Feature Extraction (8 visual signatures)**

$$\mathbf{x}_i = [E_d, L_v, \sigma_r, \kappa_r, F_h, B, C_{glcm}, H_{glcm}] \tag{1}$$

**Stage 2: Normalization**

$$\mathbf{x}_i^{\text{norm}} = \frac{\mathbf{x}_i - \boldsymbol{\mu}}{\boldsymbol{\sigma}} \tag{2}$$

**Stage 3: Graph Construction (Ultra-Chunked)**

$$G = (V, E), \quad E = \{(i, j) : j \in \text{KNN}_k(i, S)\} \tag{3}$$

**Stage 4: Clustering & Anomaly Detection**

$$\text{DBSCAN}(\epsilon = 0.8, \text{min\_samples} = 5) \tag{4}$$

**Stage 5: Classification (Random Forest on Top 5 Features)**

$$\hat{y} = \text{RandomForest}([L_v, \kappa_r, F_h, C_{glcm}, H_{glcm}]) \tag{5}$$

### B. Feature Extraction

*1) Edge Density ($E_d$):*

$$E_d = \frac{|\{p : ||\nabla I(p)|| > \tau\}|}{|P|} \tag{6}$$

Real photos have structured edges; diffusion smoothing creates different edge distributions.

*2) Laplacian Variance ($L_v$):*

$$L_v = \text{Var}(\nabla^2 I) \tag{7}$$

Measures local curvature heterogeneity. Real images: high variance; synthetic: characteristic smoothness.

*3) Noise Residual Statistics:*

$$R = I - \text{MedianFilter}(I, k = 3) \tag{8}$$

Extract standard deviation $\sigma_r$ and kurtosis $\kappa_r$. Real camera noise is Gaussian; synthetic differs.

*4) FFT High-Frequency Ratio ($F_h$) [Most Important: 34.2%]:*

$$F_h = \frac{\sum_{(u,v) \notin \mathcal{L}} |S(u,v)|}{\sum_{u,v} |S(u,v)|} \tag{9}$$

Real photos obey 1/f power law; synthetic systematically deviates. Importance: 34.2%.

*5) Blockiness Score (B):*

$$B = \frac{\sum_{i,j} w(i,j)|\nabla I(i,j)|}{\sum_{i,j} |\nabla I(i,j)|} \tag{10}$$

Detects 8×8 block boundaries from JPEG/quantization.

*6) GLCM Texture Features:*

$$\text{Contrast} = \sum_{i,j} (i-j)^2 P(i,j), \quad \text{Homogeneity} = \sum_{i,j} \frac{P(i,j)}{1 + (i - j)^2} \tag{11}$$

Real images: higher contrast; synthetic: higher homogeneity.

### C. Graph Construction with Memory Optimization

Traditional similarity matrix for 100K images requires 40 GB. We employ ultra-aggressive dual-level chunking:

$$S_{\text{full}}[i : i+c_r, j : j+c_c] = \text{cosine\_similarity}(X[i : i+c_r], X[j : j+c_c]) \tag{12}$$

Memory reduction: 40 GB $\rightarrow$ 200 MB peak = 200× improvement.

### D. Random Forest Classification

We train a Random Forest classifier on the top 5 features:

$$\text{Features} = [\text{laplacian\_var}, \text{resid\_kurtosis}, \text{glcm\_contrast}, \text{fft\_highfreq\_ratio}] \tag{13}$$

Configuration: 100 estimators, max depth 20, stratified 80/20 train/test split.

## IV. EXPERIMENTAL EVALUATION

### A. Dataset

TABLE I
DATASET COMPOSITION

| Category | Train | Test | Total |
|---|---|---|---|
| Real Images | 40,000 | 10,000 | 50,000 |
| AI-Generated | 40,000 | 10,000 | 50,000 |
| Total | 80,000 | 20,000 | 100,000 |

AI images from: Stable Diffusion (33.3%), DALL-E 3 (33.3%), Midjourney (33.3%).

### B. Results

### C. Feature Importance Analysis

**Key Finding**: Top 2 features (FFT + Kurtosis) account for 65.3% of detection performance.

### D. Clustering Analysis (100K Training Set)

22 distinct clusters identified, revealing natural image groupings.

TABLE II
RANDOM FOREST CLASSIFICATION PERFORMANCE

| Metric | Value |
|---|---|
| Accuracy | 85.65% |
| Precision | 84.48% |
| Recall | 87.34% |
| F1-Score | 0.8588 |

TABLE III
CONFUSION MATRIX (20K TEST IMAGES)

| Classification | Count | Percentage |
|---|---|---|
| True Positives (AI detected) | 8,734 | 87.34% |
| True Negatives (Real kept) | 8,395 | 83.95% |
| False Positives (Real flagged) | 1,605 | 16.05% |
| False Negatives (AI missed) | 1,266 | 12.66% |

TABLE IV
RANDOM FOREST FEATURE IMPORTANCE

| Feature | Importance | Cumulative |
|---|---|---|
| FFT High-Frequency Ratio | 34.2% | 34.2% |
| Residual Kurtosis | 31.1% | 65.3% |
| Laplacian Variance | 14.8% | 80.1% |
| GLCM Homogeneity | 10.5% | 90.6% |
| GLCM Contrast | 9.4% | 100.0% |

TABLE V
CLUSTERING RESULTS

| Metric | Value |
|---|---|
| Clusters Found | 22 |
| Noise Points | 5,132 (5.1%) |
| Clustered Points | 94,868 (94.9%) |

## V. DISCUSSION

### A. Strengths

1) **Model-Agnostic**: 85.65% accuracy across Stable Diffusion, DALL-E, Midjourney
2) **Interpretable**: Each feature has forensic meaning; importance quantified
3) **Scalable**: Processes 100K images with ¡1.2 GB peak memory
4) **High Performance**: 87.34% recall catches majority of AI images
5) **Practical**: Uses standard libraries, deployable on consumer hardware
6) **Robust**: 31.1% + 34.2% = 65.3% of decision power from just 2 features

### B. Limitations

1) **Resolution**: Features extracted at 256×256; full-resolution may differ
2) **Model Evolution**: New architectures may produce different artifacts
3) **False Positives**: 16.05% of real images incorrectly flagged
4) **False Negatives**: 12.66% of AI images missed
5) **Adversarial Robustness**: Defense-aware generation could evade detection

### C. Future Work

1) Integrate CNN embeddings from pre-trained networks for improved accuracy
2) Adversarial training for robustness against defense-aware generation
3) Extend to video analysis for temporal consistency checking
4) Multi-modal analysis combining image features + metadata + text
5) Real-time inference optimization and edge deployment
6) Comparison with other state-of-the-art detection methods

## VI. CONCLUSION

DeepTrace presents a comprehensive framework for AI-generated image detection combining unsupervised feature extraction, graph-based relationship modeling, and supervised Random Forest classification. Achieving 85.65% accuracy on 100K images with robust cross-model generalization, the system demonstrates that carefully engineered visual signatures combined with machine learning provide effective alternatives to black-box deep learning for specialized forensics applications. The ultra-aggressive memory optimization enables practical large-scale deployment. Feature importance analysis identifies FFT characteristics and noise statistics as most critical discriminators. This work demonstrates the viability of interpretable, scalable detection for emerging generative model threats.

## REFERENCES

[1] H. Farid, "Image forgery detection," *IEEE Signal Processing Magazine*, vol. 26, no. 2, pp. 16–25, 2009.
[2] S.-Y. Wang et al., "CNN-generated images are surprisingly easy to spot... for now," in *Proc. IEEE/CVF Conf. Computer Vision and Pattern Recognition*, 2020, pp. 8695–8704.
[3] R. Nightingale et al., "An introduction to convolutional neural networks," *arXiv preprint arXiv:1511.08458*, 2015.
[4] J. Zhou et al., "Graph neural networks: A review of methods and applications," *AI Open*, vol. 1, pp. 57–70, 2020.
[5] I. Goodfellow et al., "Generative adversarial nets," in *Advances in Neural Information Processing Systems*, 2014, pp. 2672–2680.
[6] A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proc. ICLR*, 2021.