

# Spectra-Temporal Attention Networks (STAN): A Dual-Stream Approach for Robust Deepfake Detection in Face Recognition Systems

Asst. Prof. Manisha Bharatram Bannagare<sup>1</sup>, Mr. Vishal Ashok Ghuge<sup>2</sup>, Mr. Rohit Vijay Shukla<sup>3</sup>,  
Mr. Kaushik Rajvilas Moon<sup>4</sup>, Mr. Om Avinash Jadhav<sup>5</sup>, Mr. Sankalp Manohar Ganvir<sup>6</sup>,  
Mr. Pratik Dnyaneshwar Shevane<sup>7</sup>, Ms. Sakshi Shankar Kewat<sup>8</sup>, Onkar Rajendra Dhanewar<sup>9</sup>

Guide, Department of Computer Science & Engineering<sup>1</sup>  
Students, Final Year Department of Information Technology<sup>2-9</sup>  
manisha180392@gmail.com

R.V. Parankar College of Engineering and Technology, Arvi, Maharashtra, India  
ghugevishal25@gmail.com and rohitvijayshukla265@gmail.com

**Abstract:** *The rapid proliferation of deep learning-based synthetic media, commonly known as "deepfakes," poses a critical threat to the integrity of biometric security systems, particularly face recognition protocols. While early generation deepfakes were easily detectable by the human eye, modern auto-encoder and diffusion-based models can generate hyper-realistic artifacts that challenge even sophisticated detection algorithms. Traditional Convolutional Neural Networks (CNNs) often fail to generalize against these threats because they over-rely on spatial pixel patterns, which are easily masked by video compression algorithms used on social media platforms. To address this limitation, this paper introduces the **Dual-Stream Spectral-Temporal Attention Network (DS-STAN)**. This novel architecture moves beyond simple pixel analysis by exploiting two fundamental weaknesses in synthetic media: the frequency-level "fingerprints" left by upsampling operations and the subtle physiological inconsistencies inherent in generated video over time. By fusing a Frequency-based stream with a Video Vision Transformer (ViT) stream, DS-STAN achieves state-of-the-art performance. Experimental results on benchmark datasets demonstrate that our model not only detects known attack types with high accuracy but also generalizes significantly better to unseen deepfake methods compared to single-modality detectors.*

**Keywords:** *Deepfake Detection, Biometric Security, Vision Transformers, Frequency Analysis, Face Anti-Spoofing, Generative Adversarial Networks*

