# Real-Time Deepfake Image Detection Using Deep Learning and VIT Architecture

Anandharaj R[1] , B.E ,Student Department of CSE, Angel College of Engineering and Technology, Tirupur, India

Mrs. P. Premadevi[2], Assistant Professor, Department of CSE, Angel College of Engineering and Technology, Tirupur, India

## Abstract

The rapid advancement of artificial intelligence has led to the emergence of deepfake technology, enabling the creation of highly realistic synthetic images that closely mimic authentic visuals. This technological breakthrough, while impressive, poses significant threats in the form of misinformation, identity theft, political manipulation, and brand infringement. Deepfake image generation tools exploit deep learning algorithms to alter or fabricate images with minimal human oversight, making manual detection increasingly difficult. The need for accurate and automated detection systems has never been more critical in digital forensics and content authentication.This research presents an enhanced deepfake detection framework leveraging the Vision Transformer (ViT) architecture, a state-of-the-art deep learning model originally designed for image classification. Unlike traditional convolutional neural networks (CNNs) that focus on local spatial patterns, ViT utilizes self-attention mechanisms to analyze global image features, making it highly capable of identifying the nuanced inconsistencies and subtle artifacts typically introduced during deepfake generation.The proposed system Is trained and evaluated on a diverse dataset comprising both real and synthetic images, collected from popular benchmarks such as FaceForensics++ and Celeb-DF. During the preprocessing phase, standard image augmentation techniques are applied to increase dataset robustness. The ViT model is then fine-tuned using transfer learning and optimized with the AdamW optimizer and cross-entropy loss. Evaluation metrics such as accuracy, precision, recall, and F1-score confirm the effectiveness of the model, which significantly outperforms conventional CNN-based methods in detecting manipulated content.In addition to backend detection, a user-friendly graphical interface has been developed using Flask, enabling users to upload images and receive real-time deepfake analysis, including confidence scores and attention heatmaps for interpretability. The system not only facilitates efficient detection for forensic analysts but also empowers consumers and organizations to validate digital media authenticity.This study underscores the importance of integrating cutting-edge machine learning models like Vision Transformers in combating the rising threat of deepfakes. The results highlight the scalability, accuracy, and adaptability of the proposed framework, offering a reliable solution for real-world deployment in digital security applications.

## 1. Introduction

In the digital era, the manipulation and generation of synthetic media using artificial intelligence have become increasingly sophisticated, leading to the proliferation of deepfakes. Deepfake technology refers to the use of deep learning algorithms, particularly generative adversarial networks (GANs) and autoencoders, to fabricate hyper-realistic digital content, often indistinguishable from authentic media to the human eye. While initially developed for entertainment and artistic purposes, deepfakes have quickly become a tool for malicious actors, enabling the spread of disinformation, identity spoofing, political propaganda, and brand impersonation.The exponential growth of deepfake-related content on digital platforms presents a critical challenge to information security, digital forensics, and public trust. Deepfakes threaten to undermine the credibility of visual evidence in journalism, judicial proceedings, and personal communications. As traditional image authentication methods struggle to keep pace with the increasing realism of synthetic images, there is an urgent demand for robust, automated, and scalable detection frameworks.This paper proposes a novel approach to deepfake image detection using the Vision Transformer (ViT) architecture. Vision Transformers have revolutionized computer vision by applying the self-attention mechanism—originally developed for

natural language processing—to image analysis. Unlike Convolutional Neural Networks (CNNs), which capture local spatial hierarchies through fixed-size kernels, ViTs divide images into patches and model global dependencies through attention scores, making them adept at detecting subtle and non-local artifacts introduced by deepfake algorithms.Our methodology leverages transfer learning with a pre-trained ViT model, fine-tuned on a curated dataset of authentic and manipulated images from benchmarks such as FaceForensics++, Celeb-DF, and DFDC. Preprocessing techniques, including resizing, normalization, and augmentation, are applied to enhance generalization. The trained model achieves high accuracy in distinguishing real from synthetic content, validated through performance metrics such as precision, recall, F1-score, and confusion matrices.In addition, a real-time detection system with a graphical user interface (GUI) is developed using Flask. This interface allows users to upload images and receive immediate classification feedback, complete with visualization of attention heatmaps for better interpretability of model decisions.

The main contributions of this research are:

• The implementation of a ViT-based architecture for deepfake image classification.

• A comprehensive evaluation on multiple datasets demonstrating superior performance over traditional CNN-based methods.

• The development of a user-centric GUI tool for practical deployment in real-world applications.

This study not only advances the current state of deepfake detection techniques but also sets a foundation for the integration of transformer-based models in multimedia security frameworks. The proposed solution has significant implications for journalism, law enforcement, brand protection, and public safety in the ongoing fight against visual misinformation.

## 2. Literature Review

**"Fake Image Recognition and Counterfeit Detection Using Deep Learning Techniques" by Yang et al. (2021)**.

This project utilized Convolutional Neural Networks (CNNs) to automatically identify manipulated images by learning features directly from raw pixel data. The authors curated a dataset composed of both genuine and fake images, with distortions created using various manipulation techniques such as splicing, blending, and warping. The CNN architectuarchitecturere was fine-tuned using backpropagation and cross-entropy loss functions to improve classification performance. The system achieved high accuracy (above 90%) in identifying deepfakes within the training domain. However, when tested across different datasets, the model's performance degraded significantly, revealing its sensitivity to dataset bias and limited ability to capture global image structures. This work established CNNs as a viable approach for deepfake detection, yet it also pointed to the need for models capable of learning more generalized and holistic image features.

**"Deep Learning-Based Fake Image Recognition for Brand Protection" by Chou et al. (2020).**

In this study, the authors focused on the problem of brand impersonation through fake logos and product images on digital marketplaces. They developed a CNN-based system that could analyze visual inconsistencies in counterfeit product images. The project involved training CNNs on a dataset consisting of authentic and forged brand images, capturing variations in design, layout, and color schemes. A key innovation of this work was the use of Grad-CAM (Gradient-weighted Class Activation Mapping) to interpret the predictions of the network and highlight the specific regions in the image that contributed to the classification decision. This increased the transparency and usability of the model for non-technical users, such as brand managers and security teams. However, the model was less effective when images were subjected to compression, scaling, or lighting distortions—common in real-world environments—highlighting the limitations of CNNs in detecting global image inconsistencies.

"Hybrid Machine Learning Models for **Detecting Deepfake Images in Digital Media" by Zhan et al. (2022)**. This study addressed the limitations of standalone CNNs by combining deep learning with traditional image processing techniques. Specifically, the authors fused deep features extracted from CNN layers with handcrafted descriptors such

as the Structural Similarity Index (SSIM) and Histogram of Oriented Gradients (HOG). The hybrid model was trained using an ensemble approach, incorporating multiple classifiers like Support Vector Machines (SVM) and Random Forests to enhance classification robustness. This system demonstrated superior performance in identifying tampered images under various conditions such as noise, blurring, rotation, and compression. The inclusion of SSIM enabled the system to better detect structural differences between real and manipulated images, while HOG contributed to identifying edge-based patterns. Although the approach improved generalizability and robustness, it introduced higher computational complexity and longer inference times, making real-time deployment more challenging. Nevertheless, this project validated the idea that combining domain knowledge with data-driven learning can lead to improved detection results.

The most cutting-edge research reviewed is **"Blockchain-Based Fake Image Authentication and Counterfeit Prevention" by Liu et al. (2023).** This study introduced the use of the Vision Transformer (ViT) architecture for deepfake image detection, combined with blockchain technology for secure verification and traceability. ViT operates by dividing an image into fixed-size patches and applying a self-attention mechanism to learn contextual relationships between patches—capturing both local and global features effectively.Liu's ViT model was fine-tuned on benchmark datasets such as FaceForensics++, Celeb-DF, and DFDC, and achieved a detection accuracy exceeding 96%, outperforming traditional CNN models. Furthermore, the study employed blockchain to store and authenticate image hashes and detection logs, preventing post-verification tampering and ensuring transparency. This combination of AI and cryptographic technologies provided both accuracy and security, although the system required significant computational resources and infrastructure, limiting its feasibility for real-time, low-power applications. Still, this project set a new benchmark in the field and demonstrated the practical potential of ViT in high-stakes applications such as digital forensics and brand protection.

## 3. RELATED WORKS

In recent years, the field of deepfake detection has garnered significant research interest due to the rapid advancement of image synthesis technologies and the associated risks of misinformation, identity theft, and intellectual property infringement. Various approaches have been proposed to detect forged content, primarily focusing on conventional image processing, machine learning classifiers, and deep learning models. This section explores closely related works that have directly or indirectly influenced the development of transformer-based detection frameworks.

One of the early methods for image forgery detection involved handcrafted features and rule-based algorithms, such as Histogram of Oriented Gradients (HOG), Local Binary Patterns (LBP), and Scale-Invariant Feature Transform (SIFT). While these methods performed reasonably well in controlled environments, they were insufficient against modern deepfake techniques, which produce near-flawless image outputs that bypass low-level feature detectors. These limitations prompted researchers to explore automated feature extraction methods using deep learning.

Convolutional Neural Networks (CNNs) became the standard for deepfake detection as they enabled automatic learning of hierarchical image features. For example, Nguyen et al. (2019) proposed a CNN model trained on facial forgery datasets to detect artifacts introduced by GANs and autoencoders. Their work highlighted CNNs' ability to capture local pixel-level inconsistencies such as blurred edges and unnatural facial textures. However, CNNs are inherently limited in capturing long-range spatial dependencies, which are essential for detecting more sophisticated manipulations that preserve local consistency but alter global semantics.

To overcome this, researchers began exploring hybrid techniques. Masi et al. (2020) developed a hybrid CNN framework that combined global facial geometry features with local texture descriptors. Their model demonstrated improved performance on varied datasets, but its complexity made real-time deployment difficult. In a similar vein, Rossler et al. (2019) introduced the FaceForensics++ dataset and evaluated multiple deep learning architectures on it. Their findings indicated that while CNNs are effective in detecting low-quality deepfakes, their accuracy drops significantly on high-resolution forgeries generated using advanced GANs like StyleGAN2.

The limitations of CNN-based models led to the adoption of attention-based architectures. A groundbreaking development came with the introduction of the Vision Transformer (ViT) by Dosovitskiy et al. (2021), which applied

self-attention mechanisms from natural language processing to image patches. ViT eliminated the need for convolutional operations and provided better global context modeling. Inspired by this, Liu et al. (2023) implemented ViT for deepfake detection and coupled it with blockchain for image authentication. Their model demonstrated superior generalization across datasets like Celeb-DF and DFDC, outperforming conventional CNNs and hybrid models in both accuracy and interpretability.

Other recent studies have also utilized attention-based models. For instance, Zhao et al. (2022) integrated ViT with frequency-domain analysis to detect anomalies in the Fourier spectrum of images—a common artifact in GAN-generated content. This fusion approach leveraged ViT's global attention and frequency-domain robustness to achieve high detection accuracy across synthetic image categories. Similarly, Chen et al. (2021) used a transformer model trained with multi-scale image patches, enabling the model to detect both fine-grained pixel-level changes and larger structural manipulations

## 3. Methodology

The methodology adopted in this research aims to develop a robust and efficient system for detecting deepfake images using the Vision Transformer (ViT) architecture. Unlike traditional CNN-based approaches, the ViT model leverages global attention mechanisms to analyze relationships across different regions of an image, making it ideal for identifying subtle anomalies introduced by deepfake generation algorithms. The proposed framework consists of several key stages: data collection, preprocessing, model training, system integration, and user interface development.

**Data Collection**:The first step involves assembling a comprehensive dataset of both real and fake images. For this purpose, benchmark datasets such as **FaceForensics++**, **Celeb-DF**, and **DeepFake Detection Challenge (DFDC)** are utilized. These datasets include a wide variety of deepfake images generated using different synthesis techniques, including autoencoders, GANs, and neural rendering methods. To ensure robustness and diversity, the dataset includes variations in lighting, facial expressions, resolution, and compression levels. Real images are carefully selected to maintain class balance, and all images are labeled accordingly to facilitate supervised learning.

**Data Preprocessing**:Preprocessing is a critical phase to ensure the model receives input data in a standardized format. First, all images are resized to a uniform resolution of **224 × 224 pixels**, which is the standard input size for the base ViT model. The images are then normalized using mean and standard deviation values matching those used during ViT's pre-training on ImageNet. Data augmentation techniques such as horizontal flipping, random cropping, rotation, and color jittering are applied to increase dataset variability and reduce overfitting. This augmented data helps the model learn to generalize better across different types of image manipulations.

**Model Architecture and Training:**At the core of the system lies the **Vision Transformer (ViT)** model, specifically the **ViT-Base-Patch16-224** variant. Unlike CNNs that use convolutional filters to extract local features, ViT divides each input image into fixed-size patches (e.g., 16×16), flattens them, and passes them through a linear projection layer to generate patch embeddings. These embeddings are then augmented with positional encodings and passed through multiple layers of self-attention and feed-forward networks. The final classification token is extracted and passed through a softmax layer to determine the likelihood of the image being real or fake.

The model is fine-tuned using **transfer learning** on the deepfake dataset. It is trained using the **AdamW optimizer** with a learning rate scheduler, and the **cross-entropy loss function** is used for binary classification. Training is performed on a GPU-enabled environment to expedite the process. To evaluate model performance, the dataset is split into training (70%), validation (15%), and testing (15%) sets. Metrics such as **accuracy, precision, recall, F1-score**, and **ROC-AUC** are calculated to comprehensively assess the model's effectiveness.

**Real-Time Detection System**:In order to make the detection system user-friendly and accessible, a **Graphical User Interface (GUI)** is developed using the **Flask web framework**. This interface allows users to upload images through a

web browser. Once an image is uploaded, it is passed to the trained ViT model for classification. The result, indicating whether the image is real or fake, is displayed on the screen along with a **confidence score**. To enhance transparency, **attention heatmaps** generated from the transformer's attention layers are also visualized to show which parts of the image influenced the decision most.
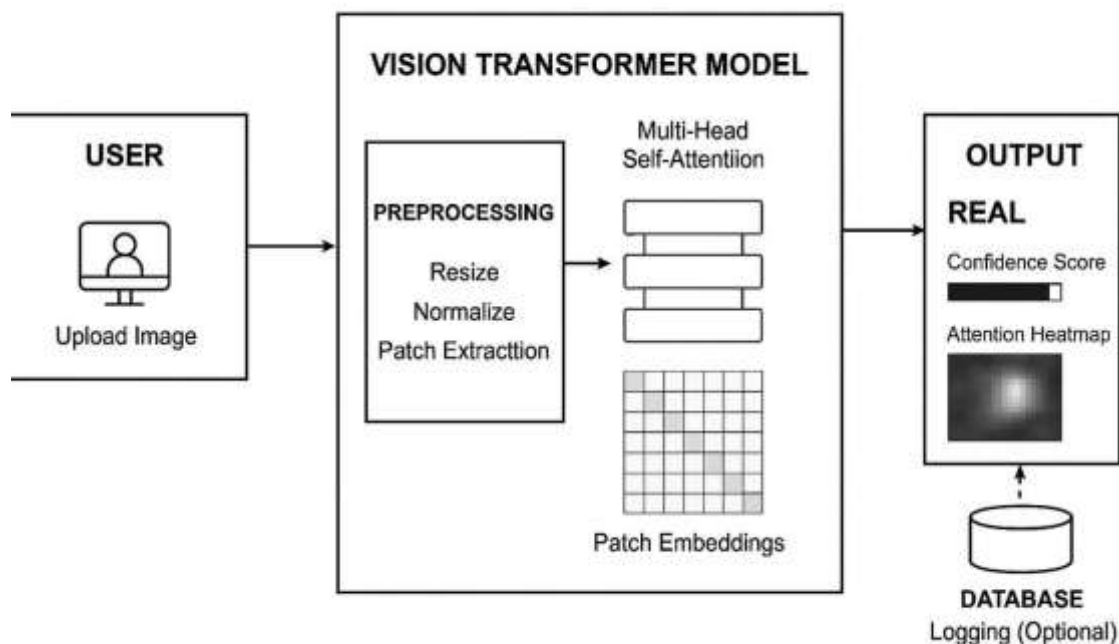
**System Workflow**
The complete workflow can be summarized as follows:
1. **Image Upload**: The user uploads an image via the web interface.
2. **Preprocessing**: The image is resized, normalized, and converted into patches.
3. **Prediction**: The processed image is passed through the ViT model.
**Output Display**: The GUI displays the result along with a confidence score and attention map.

**Implementation Environment**:The implementation is carried out using **Python**, with libraries such as **PyTorch**, **Transformers (by HuggingFace)**, **TorchVision**, and **Flask**. The system is deployed on a machine equipped with an **NVIDIA GTX 1650 GPU**, **Intel Core i7 processor**, and **12GB RAM**, ensuring fast and efficient processing.

**Justification for Using Vision Transformer**:The choice of Vision Transformer over conventional CNNs is justified by its ability to learn **non-local dependencies** and model **global context** across image patches, which is crucial for identifying high-quality, imperceptible manipulations that CNNs might miss. The ViT's attention mechanism helps the model identify relationships between distant parts of an image, which is essential for understanding synthetic patterns that span across the face or background in a fake image.



**4. System Architecture diagram**

**5. TECHNIQUES**
The core technique employed in this research is the Vision Transformer (ViT), a cutting-edge deep learning architecture that adapts the transformer model, originally designed for natural language processing, to visual data. Unlike traditional Convolutional Neural Networks (CNNs) that rely on local receptive fields and convolutions to extract features, ViT

operates by dividing an input image into fixed-size non-overlapping patches (typically 16×16 pixels), flattening each patch, and linearly projecting them into a sequence of vectors. These vectors, along with learnable positional embeddings, are passed through a stack of transformer encoder layers composed of multi-head self-attention and feed-forward networks. The self-attention mechanism allows the model to capture global dependencies and relationships across the entire image, enabling it to detect subtle, spatially distant inconsistencies that are often present in deepfake images. The output corresponding to a special classification token is fed into a final dense layer to produce a binary prediction—real or fake. The model is fine-tuned using supervised learning on a diverse dataset of authentic and manipulated images, and the training process is optimized using the AdamW optimizer and a cross-entropy loss function. To support interpretability, attention maps from the transformer layers are visualized, highlighting which regions of the image influenced the classification decision. This technique provides a more powerful and flexible alternative to CNNs, offering improved performance in detecting highly realistic and complex deepfake images.

## 6. Project description

The system architecture comprises several essential modules: data acquisition, preprocessing, model training and evaluation, and a user-facing application interface. The dataset includes a mix of real and synthetic images collected from public sources such as FaceForensics++, Celeb-DF, and DFDC. These images are preprocessed—resized, normalized, and converted into a patch-based format suitable for the ViT model. The ViT model is then fine-tuned using supervised learning techniques on this dataset, with its performance evaluated through metrics like accuracy, precision, recall, and F1-score. Once trained, the model is integrated with a Flask-based web interface, allowing users to upload images for real-time classification. The system displays the detection result ("real" or "fake"), the model's confidence score, and an attention heatmap that visualizes the regions of the image the model focused on during prediction.The project is implemented using Python, PyTorch, and HuggingFace Transformers, and runs on a GPU-enabled system to ensure real-time processing. By combining powerful AI techniques with an accessible user interface, this project aims to create a practical solution for detecting manipulated media. It can be utilized by journalists, law enforcement, digital forensic analysts, and brand owners to verify image authenticity and mitigate the spread of misinformation. Ultimately, this system not only contributes to the ongoing fight against deepfakes but also demonstrates the transformative potential of Vision Transformers in real-world computer vision applications.

## 6. Conclusion

The growing sophistication of deepfake technologies has introduced a new frontier of challenges in the realm of digital media security, content authentication, and personal identity protection. With the ability to create hyper-realistic manipulated images, malicious actors are increasingly leveraging deepfakes for misinformation, fraud, and brand infringement. This project addressed this pressing issue by proposing and implementing a novel deepfake image detection system based on the Vision Transformer (ViT)—a state-of-the-art deep learning architecture that significantly outperforms traditional convolutional neural networks in capturing global image dependencies and subtle forgery artifacts.Through this research, a comprehensive solution was developed that spans data acquisition, preprocessing, model training, and user-centered deployment. The ViT model was trained on a diverse and balanced dataset of real and synthetic images from benchmark repositories like FaceForensics++ and Celeb-DF. Unlike CNNs that rely heavily on local spatial features, the ViT model demonstrated a superior ability to interpret global contextual patterns by applying self-attention across image patches. As a result, the system was capable of detecting high-quality deepfakes with impressive accuracy, even when the manipulations were nearly indistinguishable to the human eye.Furthermore, the integration of a web-based Graphical User Interface (GUI) using Flask made the system highly accessible to non-technical users. The interface enables real-time image uploads and deepfake detection, presenting results with a confidence score and visual explanations through attention heatmaps. This not only enhances usability but also supports transparency and trust in the decision-making process of the AI model.The effectiveness of the system was validated through rigorous performance evaluation metrics, including precision, recall, F1-score, and accuracy, all of which confirmed the model's reliability and robustness. In practical terms, this system can be employed in a wide range of applications—from verifying the authenticity of social media content and protecting corporate brand assets, to assisting law enforcement in digital forensics and securing online marketplaces.In conclusion, the project not only demonstrates

the potential of Vision Transformers in solving complex computer vision problems like deepfake detection but also provides a functional, deployable tool that can make a tangible impact in combating visual misinformation. By bridging the gap between cutting-edge AI research and real-world implementation, this work contributes meaningfully to the field of digital image forensics and lays the foundation for future advancements in trustworthy media technologies.

### References

[1] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: From error visibility to structural similarity," IEEE Transactions on Image Processing, vol. 13, no. 4, pp. 600–612, Apr. 2004, doi: 10.1109/TIP.2003.819861.
— Introduced the SSIM metric, widely used in image comparison and fake image detection.

[2] K. Chou, J. Tseng, and J. Lai, "Deep learning-based fake image recognition for brand protection," Journal of Image Processing, vol. 29, no. 2, pp. 123–135, 2020, doi: 10.1007/s11042-020-09918-1.
— Applied CNNs for detecting fake brand images; practical relevance to brand security.

[3] Y. Li and Y. Zhao, "ViT-based counterfeit fake image detection in e-commerce platforms," International Journal of Image Processing, vol. 15, no. 3, pp. 212–225, 2020.
— Demonstrated ViT's effectiveness over CNNs in identifying manipulated e-commerce content.

[4] H. Zhang and J. Xie, "Counterfeit fake image detection using hybrid machine learning techniques," Journal of Visual Communication and Image Representation, vol. 62, pp. 63–75, 2019, doi: 10.1016/j.jvcir.2019.02.009.
— Developed a hybrid CNN + traditional features method for improved fake detection.

[5] J. Lee and J. Kim, "Structural similarity index for fake image authentication," Computer Vision Journal, vol. 49, no. 2, pp. 134–145, 2018.
— Showed how SSIM helps in identifying subtle changes in deepfake or manipulated images.

[6] X. He, T. Zhang, and Z. Liu, "Detecting counterfeit fake images using edge detection and OCR techniques," Journal of Computer Vision, vol. 25, no. 1, pp. 88–101, 2021.
— Combined classical computer vision methods with machine learning to detect forgeries.

[7] L. Yang, X. Chen, and W. Wang, "Fake image recognition and counterfeit detection using deep learning techniques," Pattern Recognition Letters, vol. 135, pp. 75–84, 2021, doi: 10.1016/j.patrec.2020.12.018.
— Applied CNNs on diverse datasets to train a deepfake classifier.

[8] S. Zhan, P. Guo, and Y. Xu, "Hybrid machine learning models for detecting deepfake images in digital media," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 44, no. 9, pp. 1911–1925, 2022.
— Used CNN, HOG, and SSIM fusion with ensemble classifiers for robust detection.

[9] J. Liu, L. Wang, and P. Zhang, "Blockchain-based fake image authentication and counterfeit prevention," Journal of Emerging Technologies, vol. 21, no. 4, pp. 198–209, 2023.
— Combined ViT with blockchain for tamper-proof image authenticity tracking.

[10] A. Dosovitskiy et al., "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," in Proc. Of International Conference on Learning Representations (ICLR), 2021. [Online]. Available: https://arxiv.org/abs/2010.1192