

Veriface AI: A Deep Learning-Based Approach to Deepfake Detection for Image Authenticity Validation

V.Dinesh Kumar

dept of Artificial Intelligence & Data Science

Koneru Lakshmaiah Education Foundation Vaddeswaram,India 2100080175ai.ds@gmail.com

A.kalam

dept of Artificial Intelligence & Data Science

Koneru Lakshmaiah Education Foundation Vaddeswaram,India 2100080176ai.ds@gmail.com

K.Kishore Kumar

dept of Artificial Intelligence & Data Science

Koneru Lakshmaiah Education Foundation

Vaddeswaram,India 2100080178ai.ds@gmail.com



<https://doi.org/10.55041/ijstmt.v2i2.050>

Cite this Article: Kumar, V. K. ., A. ., K. (2026). Veriface AI: A Deep Learning-Based Approach to Deepfake Detection for Image Authenticity Validation. *International Journal of Science, Strategic Management and Technology*, *Volume 10*(01). <https://doi.org/10.55041/ijstmt.v2i2.050>

License:  This article is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting use, distribution, and reproduction in any medium, provided the original author(s) and source are properly credited.

Abstract— This paper presents VeriFace AI, a deep learning-based system for detecting AI-generated fake faces. Leveraging a convolutional neural network (CNN) architecture trained on 1,400 images, the system achieves an accuracy of 83.2%. It addresses the growing concern of deepfake technology and its potential impact on digital media authenticity. We highlight the system's deployment as a web-based tool for real-time image analysis, aiming to mitigate the challenges posed by deepfake proliferation.

I. Introduction

Deepfake technology has revolutionized media synthesis, enabling the creation of highly realistic but fake images and videos. While it holds potential for creative industries, its misuse for spreading misinformation and personal exploitation has raised significant concerns. This research aims to develop an automated deepfake detection system, emphasizing accuracy, scalability, and user accessibility.

1.1 Background

The proliferation of deepfake technology has created unprecedented challenges in maintaining digital media integrity. AI-generated fake faces have become increasingly sophisticated, making manual detection nearly impossible. This research addresses the urgent need for automated detection systems. Problem Statement

The primary challenges addressed in this research include:

- Distinguishing subtle visual artifacts in AI-generated images
- Creating a robust, scalable detection system
- Developing an accessible interface for non-technical users

1.2 Research Objectives

- Design and implement a CNN-based detection model
- Achieve accuracy exceeding 80% on test data
- Create a user-friendly web interface for real-time analysis
- Evaluate system performance across diverse image types

II. Literature Review

The evolution of Generative Adversarial Networks (GANs), such as StyleGAN, has made detecting deepfakes increasingly challenging. Prior methods, including pixel-level and frequency domain analysis,

lack robustness against high-quality deepfakes. This paper builds upon these methods, employing a CNN-

Layer (type)	Output Shape	Params
conv2d (Conv2D)	(None, 126, 126, 32)	896
max_pooling2d (MaxPooling2D)	(None, 63, 63, 32)	0
conv2d_1 (Conv2D)	(None, 61, 61, 64)	18,496
max_pooling2d_1 (MaxPooling2D)	(None, 30, 30, 64)	0
conv2d_2 (Conv2D)	(None, 28, 28, 64)	36,928
flatten (Flatten)	(None, 50176)	0
dense (Dense)	(None, 64)	3,211,328
dropout (Dropout)	(None, 64)	0
dense_1 (Dense)	(None, 1)	65
Total params: 3,267,713		
Trainable params: 3,267,713		
Non-trainable params: 0		

based approach for improved accuracy.

2.1 Deepfake Technology

Recent advances in generative adversarial networks (GANs) have enabled the creation of highly realistic fake images. Studies show a 900% increase in deepfake content between 2019 and 2023.

2.2 Detection Methods

2.3 Current Limitations

Existing solutions face challenges in:

- Real-time processing
- Accuracy on high-quality fakes
- Accessibility to non-technical users

III. Methodology

The proposed system architecture consists of three layers:

1. A frontend interface for user interaction.
2. A Flask-based backend for real-time processing.
3. A CNN model for classifying real and fake faces. The model was trained on a dataset of 2,000 images (1,400 for training and 600 for testing) with balanced classes.

3.1 System Architecture

The system employs a three-tier architecture:

1. Frontend web interface
2. Flask-based backend server
3. CNN model for classification

3.2 System Architecture

- Total images: 2,000
- Training set: 1,400 images (700 per class)
- Testing set: 600 images (300 per class)
- Image resolution: 128x128 pixels

3.3 Model Architecture

IV. Implementation

4.1 Development Stack

- Frontend: HTML5, CSS3, JavaScript
- Backend: Flask (Python)
- Deep Learning: TensorFlow 2.15.0
- Image Processing: OpenCV 4.9.0

4.2 Training Process

- Optimizer: Adam
- Loss Function: Binary Cross-entropy
- Epochs: 20
- Batch Size: 32
- Learning Rate: 0.001

4.3 Deployment

The system is deployed as a web application with:

- RESTful API endpoints
- Real-time image processing
- Asynchronous prediction handling

V. Results and Discussion

The model achieved an accuracy of 83.2%, with precision, recall, and F1-score values of 0.85, 0.81, and 0.83 respectively. While effective, challenges

remain in detecting heavily edited real photos and processing high-resolution images. Future enhancements will focus on attention mechanisms and video analysis.

5.1 Performance Metrics

- Accuracy: 83.2%
- Precision: 0.85
- Recall: 0.81
- F1-Score: 0.83

5.2 Analysis

The model shows strong performance in detecting:

- GAN-generated faces
- Style-transfer modifications
- Face-swap deepfakes

5.2 Limitations

Current limitations include: Reduced accuracy on low-resolution images

- Processing time for high-resolution images
- False positives on heavily edited real photos

VI. Results and Discussion

The model achieved an accuracy of 83.2%, with precision, recall, and F1-score values of 0.85, 0.81, and 0.83 respectively. While effective, challenges remain in detecting heavily edited real photos and processing high-resolution images. Future enhancements will focus on attention mechanisms and video analysis.

6.1 Future Work

- Implementation of attention mechanisms
- Integration of temporal analysis for video
- Enhanced preprocessing pipeline
- Mobile application development

6.2 Research Directions

- Exploration of transformer architectures
- Investigation of adversarial training
- Development of explainable AI components

VII. Conclusion

VeriFace AI demonstrates the feasibility of automated deepfake detection using deep learning. The achieved accuracy of 83.2% represents a significant step toward combating digital misinformation. The web-based implementation provides a practical tool for users to verify image authenticity. VeriFace AI exemplifies the potential of CNNs in combating the deepfake crisis. By combining robust detection capabilities with user-friendly interfaces, this system provides a significant step forward in ensuring digital media integrity.

REFERENCES

- [1] Goodfellow, I., et al. (2014). Generative Adversarial Nets. NIPS.
- [2] Zhang, X., et al. (2019). DeepFake Detection Using Deep Learning.
- [3] Wang, S., et al. (2020). CNN-Based Image Analysis for Fake Detection.
- [4] Li, Y., et al. (2023). Recent Advances in DeepFake Detection.
- [5] Anderson, H., et al. (2022). Web-Based Tools for Digital Forensics.