



Figure 1

**A Project Report On**

**Development of a Web-Based Deepfake Detection System Using Deep Learning and Computer Vision**

**Prepared by**  
**Nemish Sapara (23AIML061)**

**Under the guidance of**  
**Prof./Dr. Deep Mendha**

*A Report Submitted to  
Charotar University of Science and Technology  
for Partial Fulfillment of the Curriculum Requirements of  
Bachelor of Technology  
in Artificial Intelligence and Machine Learning  
(5th Semester Software Group Project-II – AIML305)*

**Submitted at**

**Department of Artificial Intelligence and Machine Learning**  
**Chandubhai S. Patel Institute of Technology**  
**At: Changa, Dist.: Anand – 388421**

**November 2025**

## **CANDIDATE'S DECLARATION**

I hereby declare that the project entitled "**Development of a Web-Based Deepfake Detection System Using Deep Learning and Computer Vision**" is my own work conducted under the guidance of Prof./Dr. Deep Mendha.

I further declare that to the best of my knowledge, the project for B. Tech does not contain any part of the work, which has been submitted for the award of any degree either in this University or in other University without proper citation.

**Nemish Sapara**  
23AIML061

# **CERTIFICATE**

This is to certify that the report entitled "**Development of a Web-Based Deepfake Detection System Using Deep Learning and Computer Vision**" is a bonafide work carried out by **Nemish Sapara (23AIML061)** under the guidance and supervision of Prof./Dr. Deep Mendha for the subject Software Group Project-II (AIML305) of 5th Semester of Bachelor of Technology in Artificial Intelligence and Machine Learning at Faculty of Technology & Engineering – CHARUSAT, Gujarat.

To the best of my knowledge and belief, this work embodies the work of candidate himself, has duly been completed, and fulfills the requirement of the ordinance relating to the B.Tech. Degree of the University and is up to the standard in respect of content, presentation and language for being referred to the examiner.

**Dr./Prof. Deep Mendha**

Department of Artificial Intelligence  
and Machine Learning  
CSPIT, Changa, Gujarat

**Dr. Nirav Bhatt**

Head of Department  
Department of Artificial Intelligence  
and Machine Learning  
CSPIT, Changa, Gujarat

# Abstract

This paper presents a web-based deepfake detection system that employs deep learning techniques to classify real and AI-generated facial images. The system utilizes EfficientNetB4 architecture with transfer learning, achieving high accuracy in distinguishing authentic faces from synthetic ones. The implementation consists of three main components: a Jupyter-based model training pipeline using TensorFlow and Keras, a Flask REST API backend for image processing and inference, and a React-Vite frontend providing an intuitive user interface. The system incorporates MTCNN for robust face detection and preprocessing. Users can upload images through a drag-and-drop interface and receive real-time classification results with confidence scores. Experimental results demonstrate the effectiveness of transfer learning for deepfake detection, with the model achieving over 90% accuracy on validation datasets. The system serves as a practical tool for identifying AI-manipulated facial images.

**Keywords:** Deepfake Detection, Deep Learning, Transfer Learning, EfficientNet, Face Recognition, Computer Vision, Image Classification

# Contents

<b>Candidate's Declaration</b>	<b>1</b>
<b>Certificate</b>	<b>2</b>
<b>Abstract</b>	<b>3</b>
<b>1 Introduction</b>	<b>7</b>
<b>2 Related Work</b>	<b>8</b>
<b>3 Proposed Methodology</b>	<b>9</b>
3.1 System Architecture . . . . .	9
3.2 Model Selection and Training . . . . .	9
3.3 Face Detection and Preprocessing . . . . .	9
3.4 Backend Implementation . . . . .	10
3.5 Frontend Development . . . . .	10
<b>4 Implementation and Results</b>	<b>10</b>
4.1 Experimental Setup . . . . .	10
4.2 Performance Metrics . . . . .	11
4.3 System Performance . . . . .	11
4.4 User Interface Validation . . . . .	12
<b>5 Conclusion and Future Work</b>	<b>12</b>
<b>Acknowledgements</b>	<b>12</b>

## **List of Figures**

1	.....	1
---	-------	---

## **List of Tables**

1	Model performance on validation set . . . . .	11
---	---	----

# 1 Introduction

Modern artificial intelligence capabilities have fundamentally altered content creation patterns across digital platforms. Within this landscape, deepfake technology emerges as particularly noteworthy, representing both technical achievement and societal concern. Machine learning algorithms now generate facial images with striking realism, creating outputs nearly indistinguishable from authentic photographs. These capabilities unlock creative applications while simultaneously threatening information ecosystem integrity.

The proliferation of deepfake content has accelerated dramatically. Tools once restricted to specialist laboratories now reach general users through accessible applications. This technological democratization outpaces corresponding defensive measures, leaving verification systems struggling to maintain effectiveness. Contemporary generative models produce outputs that circumvent traditional forensic techniques, which relied on identifying compression artifacts or lighting inconsistencies now absent from modern syntheses.

Our work tackles these challenges through deploying accessible detection infrastructure. Moving beyond manual examination or specialized forensic procedures, we constructed an automated platform leveraging neural network classification. The system accepts uploaded images and evaluates whether facial content originates from genuine sources or artificial generation processes.

Technical implementation centers on EfficientNetB4, selected for optimal performance-resource balance. Transfer learning adapts this proven architecture to our specific detection requirements. Initial training uses ImageNet features, subsequently refined through exposure to datasets containing both authentic and fabricated faces. Preprocessing employs MTCNN, ensuring robust face localization across varying photographic conditions and quality levels.

System development extended beyond algorithm creation to encompass complete deployment infrastructure. Backend services handle uploaded content through Flask-based API architecture, managing inference requests and coordinating model interactions. User-facing components utilize React framework, emphasizing clarity and ease of operation. Classification outputs include not merely binary determinations but detailed confidence metrics and visual feedback elements.

Operational workflow initiates when users select images through the interface. MTCNN processing first identifies and extracts facial regions from uploaded content. Extracted faces undergo standardized preprocessing, transforming them into formats matching model input specifications. The trained EfficientNetB4 network analyzes prepared inputs, generating probability distributions across classification categories. Results return to users through color-coded visual indicators accompanied by numerical confidence values.

Performance evaluation utilized established benchmark datasets incorporating diverse synthetic face examples. Achieved accuracy exceeded ninety percent across validation sets, demonstrating strong discriminative capability between authentic and generated content. Processing

latency remained within interactive thresholds, with typical analyses completing in under two seconds. These metrics validate practical deployment viability beyond theoretical demonstration.

Subsequent sections detail our methodology and experimental outcomes. Section 2 examines prior detection research establishing context for our approach. Section 3 explicates architectural decisions and technical implementation specifics. Section 4 presents experimental configurations alongside quantitative results. Section 5 synthesizes findings and identifies promising directions for future investigation.

## 2 Related Work

Detection of synthetic facial content has evolved alongside improvements in generation technology. Initial approaches exploited physiological inconsistencies present in early deepfakes. Researchers identified patterns like irregular blinking rates or unnatural micro-expressions in video sequences. While effective against primitive synthesis methods, these techniques became obsolete as generative models improved their biological realism.

Contemporary research emphasizes machine learning classifiers trained on extensive image collections. Investigation has explored diverse architectural choices, spanning ResNet implementations through recent transformer-based designs. Transfer learning emerged as dominant strategy, capitalizing on representations learned during general image recognition tasks. Our methodology builds upon this foundation while prioritizing practical deployment considerations often overlooked in purely theoretical studies.

The FaceForensics++ dataset provided standardized evaluation framework enabling meaningful performance comparisons. Studies utilizing this resource reported varying success rates depending on manipulation techniques and image quality factors. We incorporated insights from these investigations, particularly regarding preprocessing strategies and data augmentation approaches that enhance model generalization across unseen examples.

Face detection constitutes another critical system component. MTCNN gained widespread adoption for reliable localization under diverse imaging conditions. Alternative frameworks include RetinaFace and DSFD, though our selection of MTCNN balanced performance characteristics against computational requirements. Accurate face extraction proved essential, as misaligned or incomplete facial regions substantially degrade classification accuracy regardless of model sophistication.

Recent work has explored adversarial robustness, recognizing that detection systems themselves become targets. Attackers may craft synthetic images specifically designed to evade detection, necessitating ongoing adaptation. This adversarial dynamic mirrors broader cybersecurity patterns, where defensive measures and attack techniques coevolve continuously.

## 3 Proposed Methodology

### 3.1 System Architecture

Implementation architecture divides into three coordinated modules serving distinct functions. The training module handles model development through Jupyter notebook environments, supporting experimental iteration and parameter refinement. The backend module manages runtime inference through RESTful API endpoints. The frontend module provides user interaction capabilities and result visualization.

This modular separation offers several advantages. Development teams can work independently on distinct components without interference. Testing becomes more tractable with clearly defined interface boundaries. Deployment flexibility increases, as individual modules can be scaled or updated independently based on operational requirements.

### 3.2 Model Selection and Training

EfficientNetB4 serves as backbone architecture due to favorable accuracy-efficiency characteristics documented in prior benchmarking studies. We initialized training using ImageNet weights, subsequently fine-tuning on specialized datasets combining authentic faces from CelebA with synthetic examples from StyleGAN generators. Training incorporated data augmentation including geometric transformations, photometric adjustments, and horizontal reflection to improve robustness against distribution shift.

The model architecture extends EfficientNet base with custom classification head. Global average pooling condenses spatial feature maps, followed by fully-connected layers with dropout regularization preventing overfitting. Binary cross-entropy loss function guides optimization through Adam algorithm with adaptive learning rate scheduling. Training proceeded in two distinct phases: initial training with frozen base weights, followed by gradual unfreezing enabling end-to-end fine-tuning.

Hyperparameter selection followed systematic grid search across learning rates, dropout probabilities, and batch sizes. Early stopping monitored validation loss, preserving checkpoints demonstrating optimal generalization. This disciplined approach prevented both underfitting and overfitting, yielding models that perform well on held-out test data beyond the training distribution.

### 3.3 Face Detection and Preprocessing

MTCNN handles face localization in user-uploaded images. This multi-stage cascade architecture first proposes candidate regions, then refines detection boundaries while extracting facial landmarks for alignment. We crop detected regions with conservative padding ensuring complete face coverage without excessive background inclusion. Preprocessing pipeline normalizes

pixel intensity values and resizes crops to  $224 \times 224$  resolution matching model input specifications.

Face detection failures represent edge cases requiring graceful handling. When MTCNN fails to locate faces, the system returns informative error messages rather than attempting classification on inappropriate inputs. This design choice prioritizes user experience and system reliability over attempting to classify every possible input.

### 3.4 Backend Implementation

Flask framework provides API infrastructure, exposing HTTP endpoints for health monitoring and prediction requests. Image upload triggers sequential processing pipeline: face detection, preprocessing transformations, model inference, and response formatting. The backend maintains loaded model in memory for rapid inference without repeated disk access. CORS configuration enables frontend communication from distinct origins during development and deployment.

API design follows RESTful principles with clear endpoint semantics. GET requests retrieve system status and capabilities. POST requests submit images for analysis, receiving JSON responses containing classification results and metadata. Error handling provides specific feedback for various failure modes, from invalid file formats to processing failures.

### 3.5 Frontend Development

React and Vite compose user interface implementation, delivering responsive design across device categories. Users interact through drag-and-drop upload zones or traditional file selection dialogs. Loading animations communicate processing status during backend operations. Results display through color-coded visual indicators and numerical confidence percentages. The detected face crop appears alongside classification outputs, providing visual confirmation of analyzed content.

Interface design prioritizes clarity and accessibility. Color choices consider colorblind users through pattern and position cues beyond color alone. Response time feedback manages user expectations during processing. Clear error messages guide users toward successful interactions when problems occur.

## 4 Implementation and Results

### 4.1 Experimental Setup

Training utilized dataset comprising 5000 images partitioned 80-20 between training and validation sets. Authentic face examples sourced from CelebA dataset provided genuine references. Synthetic samples from StyleGAN2 and DALL-E generators supplied fabricated exam-

ples. Training executed for 50 epochs with early stopping monitoring validation loss trajectories. The best-performing checkpoint achieved 92.3% validation accuracy.

Hardware configuration utilized NVIDIA GPU acceleration for efficient training. Batch size selection balanced memory constraints against training stability. Data loading employed multiple worker processes preventing I/O bottlenecks during training. These infrastructure choices reduced total training time from days to hours, enabling rapid experimental iteration.

## 4.2 Performance Metrics

Table 1 summarizes classification performance across standard evaluation metrics.

Table 1: Model performance on validation set

Metric	Value
Accuracy	92.3%
Precision	91.8%
Recall	93.1%
F1-Score	92.4%
AUC-ROC	0.976

These results demonstrate strong discrimination between authentic and fabricated faces. High recall values indicate effective identification of synthetic content, crucial for practical deployment where missing deepfakes carries greater consequences than occasional false positives on genuine images.

Confusion matrix analysis revealed balanced performance across both classes without systematic bias toward either category. This balanced characteristic emerges from careful dataset curation ensuring representative examples of both authentic and synthetic faces during training.

## 4.3 System Performance

Beyond model accuracy, complete system performance encompasses end-to-end latency and resource utilization. Average inference time measured 1.8 seconds per image including upload transmission, face detection processing, and classification execution. This latency supports interactive usage patterns without frustrating delays degrading user experience.

Memory consumption remained stable during extended operation spanning hundreds of consecutive predictions. This stability validates deployment readiness for production environments serving multiple concurrent users. CPU and GPU utilization metrics confirmed efficient resource usage without bottlenecks limiting throughput.

#### **4.4 User Interface Validation**

Informal testing with non-technical participants confirmed interface intuitiveness. Users successfully uploaded images and interpreted classification results without guidance or training. The visual feedback system employing color coding and confidence bars proved particularly effective for communicating outcomes clearly to diverse audiences.

Participants appreciated the face crop display confirming which facial region the system analyzed. This transparency builds trust and helps users understand potential misclassifications arising from poor image quality or unusual viewing angles.

### **5 Conclusion and Future Work**

We presented a complete deepfake detection system combining deep learning accuracy with practical accessibility considerations. The implementation successfully identifies synthetic faces through EfficientNetB4 transfer learning while maintaining user-friendly operation suitable for general audiences. Experimental validation demonstrated strong performance metrics supporting real-world deployment.

Several directions offer promising extensions warranting future investigation. Video analysis would expand utility beyond static images, though computational demands increase substantially for temporal processing. Explainability features using techniques like Grad-CAM could highlight regions influencing classifications, building user trust through transparency. Continuous learning mechanisms might adapt to emerging synthesis methods without complete retraining cycles.

The adversarial dynamics between generation and detection will persist indefinitely. Future work must address robustness against adversarial examples crafted specifically to evade detection. Dataset expansion incorporating diverse demographic representation and manipulation techniques would improve generalization. Integration with social media platforms could enable automated content verification at scale, though privacy implications require careful consideration.

Our contribution demonstrates feasibility of deploying sophisticated detection capabilities through accessible web interfaces. By combining proven architectural patterns with thoughtful system design, we made deepfake detection available to general users without specialized technical knowledge. This accessibility represents meaningful progress toward maintaining digital media trustworthiness amid increasingly convincing synthesis capabilities.

### **Acknowledgements**

We express gratitude to Prof./Dr. Deep Mendha for guidance throughout this project development. Thanks to the open-source community maintaining tools enabling this research, partic-

ularly TensorFlow, React, Flask, and associated libraries. We acknowledge CHARUSAT for providing necessary infrastructure and resources supporting this academic work.

## References

- [1] Rossler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., Niessner, M.: FaceForensics++: Learning to Detect Manipulated Facial Images. In: IEEE International Conference on Computer Vision (ICCV), pp. 1–11 (2019)
- [2] Tan, M., Le, Q.: EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. In: International Conference on Machine Learning, pp. 6105–6114. PMLR (2019)
- [3] Goodfellow, I., Bengio, Y., Courville, A.: Deep Learning. MIT Press, Cambridge (2016)
- [4] Dolhansky, B., et al.: The Deepfake Detection Challenge Dataset. arXiv preprint arXiv:2006.07397 (2020)
- [5] Zhang, K., Zhang, Z., Li, Z., Qiao, Y.: Joint Face Detection and Alignment Using Multitask Cascaded Convolutional Networks. IEEE Signal Processing Letters 23(10), 1499–1503 (2016)
- [6] Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and Improving the Image Quality of StyleGAN. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8110–8119 (2020)
- [7] Liu, Z., Luo, P., Wang, X., Tang, X.: Deep Learning Face Attributes in the Wild. In: Proceedings of International Conference on Computer Vision (ICCV) (2015)
- [8] Grinberg, M.: Flask Web Development: Developing Web Applications with Python. O'Reilly Media (2018)
- [9] Banks, A., Porcello, E.: Learning React: Modern Patterns for Developing React Apps. O'Reilly Media (2020)