



# **Image Deepfake Detection Through Machine Learning–Driven Visual Analysis**

Final Project Report Submitted to  
The Department of Computer Science  
Faculty of Computer and Information Technology  
Jordan University of Science and Technology  
In Partial Fulfillment of the Requirements for the Degree of Bachelors of Science in Computer Science

Prepared by:

Ahmad Kamhia[143397]  
Ameen Jarrar [161811]  
Yousef Mohammad Alsaïd [158304]  
Awn Ali Taani [157854]

Supervisor:  
Farah Alshanik

January 2026

# Image Deepfake Detection Through Machine Learning–Driven Visual Analysis

\*Note: Sub-titles are not captured in Xplore and should not be used

1<sup>st</sup> Ahmad Kamhia

Department of Computer Engineering  
Jordan University of Science and Technology  
Irbid, 22110, P. O. Box. 3030  
qaahmad22@cit.just.edu.jo

2<sup>nd</sup> Ameen Jarrar

Department of Computer Engineering  
Jordan University of Science and Technology  
Irbid, 22110, P. O. Box. 3030  
amjarrar22@cit.just.edu.jo

3<sup>rd</sup> Yousef Alsaid

Department of Computer Engineering  
Jordan University of Science and Technology  
Irbid, 22110, P. O. Box. 3030  
ymalsaid210@cit.just.edu.jo

4<sup>th</sup> Awn Taani

Department of Computer Engineering  
Jordan University of Science and Technology  
Irbid, 22110, P. O. Box. 3030  
aataani21@cit.just.edu.jo

**Abstract**—Man-made, AI-produced deepfakes are becoming increasingly realistic, more available, and they raise many questions and concerns, such as: misinformation, fraud, and identity theft, etc. Many deepfake detection systems have been developed, but most tests are performed on ideal, clean datasets, which are not representative of images found on these platforms (e.g., Instagram, Twitter, WhatsApp). Images that are uploaded on social media platforms often have some type of compression and are rescaled, or slightly degraded, which ultimately reduce the performance of "standard" deepfake detection techniques. Thus, the development of better methods for the identification of manipulated images, especially in real-world environments, creates an immediate need to address this problem, which is especially important for IT Security, Digital Forensics, and the Verification of Content Online.

This body of research addresses only image-based deepfake detection. We use publicly available datasets, namely Celeb-DF v2 and FaceForensics++, and we augment the datasets to create simulated conditions that replicate how images appear when printed on social media platforms, including JPEG compression, resizing of images, and noise/blur added to the images. We developed a deep learning framework using convolutional neural networks (CNNs) in conjunction with transfer learning (using a pre-trained model) to improve feature extraction. The model was trained on our augmented dataset and tested against new images to determine how effectively and accurately it could generalise to recognise manipulated images across a variety of degradations when critical factors are considered.

The experimental findings confirm that this proposed model provides accurate detection of deepfake images with or without compression, size, or added noise. When evaluated against a baseline or other forms of transfer learning, the proposed method exhibited superior strength in all simulated social media environments. The model's use of attention maps and visualizations confirms that the model can determine the changed areas of an image, making it a reliable and practical tool for verifying deepfakes in actual practice.

This thesis provides evidence that a comprehensive image-based approach to deepfake detection can be utilized for analyzing real-world, social media-like images. Through

enhancing pre-existing datasets and utilizing current CNNs combined with transfer learning, a superiorly proficient and dependable deepfake detection methodology is demonstrated. Included in this research is a collection of professional quality dataset resources, a strong detection methodology, and recommendations for future use of these results for IT security and digital forensics. The findings of this research support the need for investigations of possible image degradation to be incorporated in the design and development of future deepfake detection systems, and a framework for the enhancement of artificial intelligence (AI) techniques to authenticate digital content.

**Index Terms**—Deepfake detection, Image-based detection, Convolutional neural networks (CNNs), Transfer learning, Social media images, Noise and blur augmentation, Robust detection, Digital forensics

## I. PROJECT GOALS AND OBJECTIVES

A primary objective of the current project is to build a robust and precise method to identify altered digital content, or "deepfake" images posted on the Internet, through the application of machine learning technology.

**Project goals consist of :**

- Evaluating the efficacy of existing techniques for identifying deepfake images when images are used under non-ideal circumstances (e.g. compression, resizing, noise, blurriness, etc. from social media platforms).
- Developing an augmented dataset created from currently available publicly available datasets (Celeb-DF v2 and FaceForensics++) of deepfake images, using realistic social media image deterioration (JPEG compression, scaling of images, noise, and blurriness).
- Creating a CNN (convolutional neural network) that is used as a method for detecting deepfake images, utilising transfer learning as a mechanism for improving feature extraction and generalisation.

- Assessing the fidelity and generalisability of the proposed methods compared to traditional baseline techniques as well as comparing these methods to similar alternative methods based on transfer learning.
- Utilising attention maps and visualisation methods to aid understanding of predictions made by the developed models and verify that facial areas were correctly identified.
- Showing how the blending framework designed within the current project's context has potential applications within the - IT Security, digital forensics and online verification of digital content.

## II. INTRODUCTION

The rapid progress of artificial intelligence (AI), particularly in generative models, has led to the emergence of highly realistic synthetic media known as deepfakes. Deepfakes refer to manipulated or AI-generated images and videos that convincingly imitate real individuals, making them difficult to distinguish from authentic content.

The technologies covered pose an obvious risk due to questions about misinformation, fraud, identity theft, and privacy violations. Trust in digital media is diminished by these risks [1], [2]. Furthermore, as techniques for generating deepfakes have become more advanced and accessible, there is a growing need for reliable, effective ways of identifying them. Therefore, there are a variety of methods for detecting deepfakes that have been proposed for researchers to consider developing.

Deep learning-based approaches, especially CNNs, have received considerable attention because they have demonstrated the ability to produce good results when identifying manipulated facial images [3], [4]. There is a lot of focus in current studies on improving the accuracy of these models through the development of advanced architectures for the CNNs, using a large volume of training data, and other means [5], [6].

Recent reviews of the literature demonstrate how many researchers have entered this rapidly developing field of study, and they highlight the necessity of developing robust detection systems capable of keeping pace with the emergence of new deepfake techniques [7], [8].

While it is clear that advances have been made in this area of study, the vast majority of existing detection systems have been created based on and trained from a limited amount of high-quality, pristine data. As can be seen quite clearly in practice on social media channels, images that are frequently shared on platforms such as Instagram, Twitter, and WhatsApp have often been subjected to various forms of degradation due to JPEG compression, resizing, noise, and blur.

Due to this degradation, the performance of the majority of conventional detection models has been adversely affected, significantly decreasing their overall performance in addition to their generalisability and robustness [9], [10].

Currently, there is still a significant lack of research on the creation of reliable deepfake detection tools that will operate under similar conditions to those experienced on social

media. However, recent developments using a transfer learning approach will allow for better training and robustness in the detection tools used to identify deepfakes [11].

Despite these improvements, most methods to date do not consider how the combined impact of compression, resize, and noise that occurs during the process of sharing digital content on the internet affects the end-user's ability to properly assess an image's authenticity. Furthermore, the area of interpretability and explainability has received very little focus in the area of developing a digital forensics and security-based deepfake detection model [12], [13].

This research will focus exclusively on detecting image-based deepfakes. The overall goal of this research is to create an effective bridge between laboratory testing of the model's ability and the final product's field performance. Publicly available datasets, such as Celeb-DF v2 and FaceForensics++, will be augmented with various levels of quality degradation (due to image resizing, noise, and compression) that are common to social media platforms. The augmented dataset will then be used to create a new CNN-based deep-learning framework where the majority of the layers will utilize transfer learning to improve upon feature extraction and generalizability and to provide sufficient robustness across multiple image-quality levels. A number of different techniques, including visualization of CNN outputs, will be conducted to aid in identifying which areas of the model are being trained to identify features related to manipulated facial areas. This will provide additional assurance of practical applicability for use in the areas of digital forensics and security [14], [15].

The remainder of this paper is organized as follows. Section II presents a review and analysis of related work. Section III the proposed approach and methodology, including the tools, algorithms, and overall system pipeline.

### A. REVIEW AND ANALYSIS OF RELATED WORK

Deepfake generation and detection have attracted growing research interest due to the rapid advancement of Artificial Intelligence (AI) and Deep Learning (DL) techniques. Several studies have investigated the characteristics, risks, and detection mechanisms associated with deepfake media.

The authors of [1], provided a complete overview of the generation and detection of deepfakes through the analysis of the different techniques available for the generation and detection of deepfakes, with respect to visual-based, temporal-based and frequency-based characteristics. This paper also highlighted both the strengths and weaknesses associated with the use of deep learning-based methods for the detection of deepfakes. One area where deep learning-based methods are weak is when applied to previously unseen manipulation methods/models and dataset bias that may adversely affect their accuracy. Similarly, Karasavva and Noorbhai [2] examined the ethical and social concerns surrounding deep fakes, focusing primarily on the policy challenges and potentially serious negative impacts that could result from manipulated media. There are a number of works which focus on enhancing the accuracy of detecting deepfake content via the use of convolutional

neural networks (CNNs). Patel et al. [3] proposed the use of an enhanced version of a dense CNN architecture with respect to improved reuse of features leading to improved accuracy when identifying deepfake image content. Abdullah et al. [4] have evaluated the most current developments in deepfake detection and evaluated the existing techniques that identify deepfake image data in the face of ever changing manipulation methods and demonstrated that the accuracy of existing models is diminishing as these manipulation techniques evolve. Patel et al. [5] also provided a discussion of the issues confronting deepfake generation and detection, highlighting issues related to dataset diversity and generalization.

Numerous reviews and bibliometric analyses have further demonstrated the rapid expansion of this research field. Gupta et al. [6] presented a review of the state-of-the-art techniques for deepfake detection, focusing on approaches based on fusion and advanced machine learning methodologies; Acim et al. [7] performed a bibliometric study of the past 10 years on deepfake research trends, finding that robustness and real-world applicability are key areas of difficulty in the field. Tolosana et al. [8] provided a comprehensive review of face manipulation and detection test bed datasets and underscored the need for standardised evaluation methods. Although many studies have reported promising results when evaluating the performance of deepfake detectors on test data, they have also noted that in-the-wild and real-world degradation of test data degrades deepfake detector performance considerably. Diallo et al. [9] showed that JPEG compressor artefacts served as vulnerabilities to the ability of Convolutional Neural Networks (CNN) to detect manipulated content. In a similar manner Qazi et al. [10] proposed a deep learning-based method for the detection of image forgeries that incorporated transfer learning to achieve highly accurate results on "clear" (non-degraded) datasets; social media-like degradation was not explicitly included in the proposed method. Ranjan et al. [11] established the effectiveness of transfer learning based CNN approach at increasing the ability of the model to generalise across multiple datasets; while the authors stressed that combined degradation occurred from the processes of compression, resizing, and addition of background noise, they did not model these losses together. Recent advancements in this area have emphasised the need for researchers to focus more on interpretability and human centred analysis of models.

Diel et al. [12] discuss the damaging psychological and behavioural effects of the technology behind the development of Deepfakes. Meanwhile, Abdel-Wahab and Alkhatib [13] point out that developing robust and explainable detection methods will be crucial for the successful implementation of Deepfake technology in real-life scenarios. In general, most of the current Deepfake detection techniques are capable of achieving very high levels of accuracy in a controlled environment however, their effectiveness diminishes when used with real-world degraded images. Our study aims to mitigate some of these limitations by focusing on the image-based detection of Deepfakes under typical social media conditions. We accomplish this by applying data augmentation, transfer

learning, and visualisation techniques to increase our robustness and practical usability when detecting Deepfakes under these conditions.

Overall, existing deepfake detection approaches achieve high accuracy under controlled conditions but often lack robustness and interpretability when applied to real-world, degraded images. This work addresses these limitations by focusing on image-based deepfake detection under realistic social media conditions, incorporating data augmentation, transfer learning, and visualization techniques to enhance robustness and practical applicability.

### *B. Significance of work*

The fast-growing development of Generative Artificial Intelligence (AI) means that the development of realistic, wide-reaching, and almost impossible-to-detect Deepfakes is rapidly accelerating, thereby threatening the trust, privacy, and security of the digital world. Existing Deepfake detection systems are trained to achieve a very high level of accuracy when tested against unadulterated, ideal data sets; however, the accuracy of these systems decreases significantly when the input pictures are subject to very common and well known social media modifications such as resizing, noise, compression or blurring. This means that there is currently a huge difference between pure laboratory testing and applying new methodologies in the real-world.

The primary focus of this study is to provide an alternative approach to deepfake detection that can be used in the real world. To that end, the researchers augmented publicly available datasets to create simulated versions of social media degradation. They then developed a CNN (Convolutional Neural Network)-based detection framework that uses transfer learning to achieve high levels of accuracy under realistic conditions. The study also utilises attention maps, as well as visualisation techniques, to improve the interpretability of the models developed in this study. This will be essential for these models to be implemented in topics related to digital forensics, cyber security and content verification.

The results of this research have not only improved the detection capabilities of deepfakes technically but have also provided the AI and cyber security communities with: A method for bridging the gap between laboratory and real-world testing so that detection models can be improved when evaluating under cases of realistic image degradation. A reproducible methodology that combines data augmentation, convolutional neural network detection, and transfer learning; this method can be extended or adapted to other types of image manipulation. A reliable way to verify images posted on social media, which will assist digital forensics and IT security. Encouragement to continue researching next-generation deepfake detection systems that can detect with higher accuracy and have increased interpretability and robustness as the technologies improve over time. This study provides a robust and practical methodology for detecting images based on their real-world applications and moving beyond a purely academic

Paper	Methodology	Key Results	Limitations
[1]	Comprehensive survey of deepfake generation and detection	Provides taxonomy and comparative analysis of methods	No experimental validation or robustness evaluation
[2]	Policy and threat analysis of deepfake misuse	Highlights social, legal, and ethical risks	Does not address technical detection methods
[8]	Survey of face manipulation and fake detection techniques	Benchmarks datasets and detection pipelines	Limited focus on robustness to image degradation
[3]	Improved Dense CNN for deepfake image detection	Achieves high accuracy on benchmark datasets	Trained and tested on clean images only
[5]	Case study of deepfake generation and detection challenges	Identifies performance trade-offs and dataset bias	Lacks degradation-aware evaluation
[6]	Review of ML, DL, and fusion-based detection approaches	Shows effectiveness of hybrid models	No explicit modeling of social media distortions
[11]	Transfer learning-based CNN framework	Improved cross-dataset generalization	Does not consider compression or noise effects
[9]	CNN supervision with compression-aware training	Robust detection under JPEG compression	Ignores resizing, blur, and noise
[10]	ResNet50v2 with transfer learning for image forgery detection	Achieves high accuracy on CASIA datasets	Not evaluated on deepfake-specific datasets
[13]	Review of detection and prevention techniques	Highlights emerging methods such as XAI and blockchain	No implementation or experimental validation
<b>This Work</b>	CNN with transfer learning and degradation-aware augmentation	Robust detection under compression, resizing, noise, and blur	Limited to image-based deepfakes

TABLE I

EXTENDED SUMMARY AND COMPARISON OF RELATED WORK IN IMAGE-BASED DEEPFAKE DETECTION. THIS TABLE SPANS BOTH COLUMNS IN IEEE STYLE.

approach of utilizing the "ideal" datasets to create trustworthy AI-enabled media verification platforms.

### III. APPROACH AND METHODOLOGY

#### A. Methodology

This work proposes a CNN-based deep learning framework for **image-based deepfake detection**, designed to operate robustly under social media-like degradations. The methodology involves three main components: dataset preparation, model architecture, and training strategy.

1) *Dataset Preparation*: Publicly available datasets, including **Celeb-DF v2** and **FaceForensics++**, are augmented to simulate real-world conditions encountered on social media. Augmentations include:

- **JPEG compression** to simulate lossy image uploads.
- **Resizing and scaling** to model display and platform effects.
- **Noise and blur** to represent camera artifacts and transmission degradation.

Let  $I_{\text{orig}}$  be an original image. After applying degradation  $D$ , the processed image  $I_{\text{deg}}$  is:

$$I_{\text{deg}} = D(I_{\text{orig}}) \quad (1)$$

where  $D$  can be any combination of compression, scaling, noise, or blur transformations.

2) *CNN Architecture*: A **transfer learning-based Convolutional Neural Network (CNN)** is employed to enhance feature extraction and reduce training time. The network backbone is pre-trained on ImageNet and fine-tuned using the augmented deepfake datasets. Key components include:

- **Convolutional layers** for spatial feature extraction.
- **Pooling layers** for dimensionality reduction.
- **Fully connected layers** for classification into real or fake categories.
- **Attention visualization** to analyze regions of the face that influence the model's decision.

Figure 1 illustrates the overall pipeline of dataset augmentation, CNN feature extraction, and classification.

3) *Training Strategy*: The model is trained with **cross-entropy loss**, optimized using the **Adam optimizer**, and validated using **k-fold cross-validation** to ensure robustness across diverse conditions. Performance metrics include **accuracy, precision, recall, and F1-score**.

#### B. Location and Safety Considerations

Since this study deals with **publicly available image datasets**, there are no physical safety concerns. Ethical considerations involve ensuring that all datasets comply with privacy and consent regulations. Model deployment is limited to detection purposes, avoiding misuse of generated or manipulated images.

#### C. Expected Results / Outputs

The proposed system is expected to:

- Accurately classify images as real or fake, even under degradation effects.
- Maintain robust performance across different datasets and social media-like transformations.
- Provide visual explanations for its predictions, enabling interpretability and forensic utility.

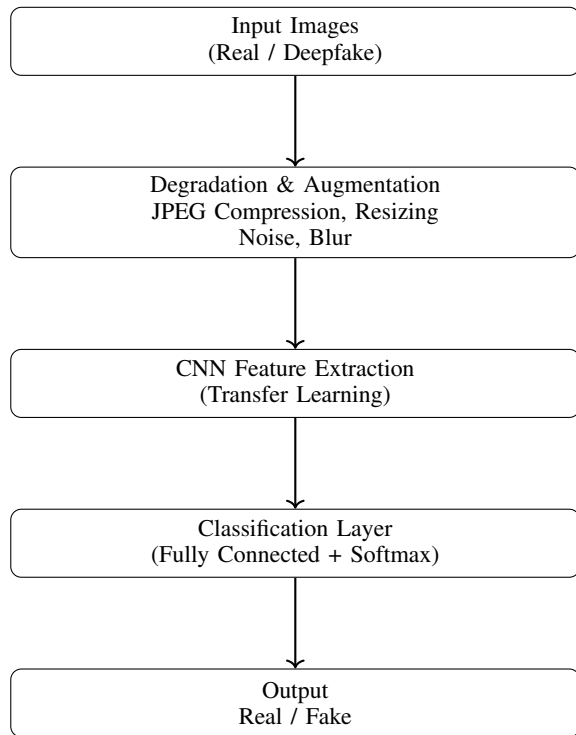


Fig. 1. Proposed deepfake detection pipeline including degradation-aware augmentation, CNN-based feature extraction, and classification.

- Serve as a benchmark for future research in robust, degradation-aware deepfake detection.

Outputs include:

- 1) Classification labels (real/fake) for each input image.
- 2) Probability scores indicating confidence levels of predictions.
- 3) Attention maps highlighting critical facial regions influencing detection.

## REFERENCES

- [1] T. Zhang, "Deepfake generation and detection: A survey," *Multimedia Tools and Applications*, vol. 81, pp. 6259–6276, Feb. 2022.
- [2] V. Karasavva and A. Noorbhai, "The real threat of deepfake pornography: A review of canadian policy," *Cyberpsychology, Behavior, and Social Networking*, vol. 24, no. 3, pp. 203–209, 2021.
- [3] Y. Patel, S. Tanwar, P. Bhattacharya, R. Gupta, T. Alsuwian, I. E. Davidson, and T. F. Mazibuko, "An improved dense cnn architecture for deepfake image detection," *IEEE Access*, vol. 11, pp. 22081–22095, 2023.
- [4] S. M. Abdullah *et al.*, "An analysis of recent advances in deepfake image detection in an evolving threat landscape," in *Proceedings of the IEEE Symposium on Security and Privacy (SP)*, (San Francisco, CA, USA), pp. 91–109, 2024.
- [5] Y. Patel *et al.*, "Deepfake generation and detection: Case study and challenges," *IEEE Access*, vol. 11, pp. 143296–143323, 2023.
- [6] G. Gupta, K. Raja, M. Gupta, T. Jan, S. T. Whiteside, and M. Prasad, "A comprehensive review of deepfake detection using advanced machine learning and fusion methods," *Electronics*, vol. 13, no. 1, 2024.
- [7] B. Acim, M. Boukhelif, H. Ouhnni, N. Kharmoum, and S. Ziti, "A decade of deepfake research in the generative ai era, 2014–2024: A bibliometric analysis," *Publications*, vol. 13, no. 4, 2025.
- [8] R. Tolosana, R. Vera-Rodriguez, J. Fierrez, A. Morales, and J. Ortega-Garcia, "Deepfakes and beyond: A survey of face manipulation and fake detection," *Information Fusion*, vol. 64, pp. 131–148, 2020.

- [9] B. Diallo, T. Urruty, P. Bourdon, and C. Fernandez-Maloigne, "Robust forgery detection for compressed images using cnn supervision," *Forensic Science International: Reports*, vol. 2, p. 100112, 2020.
- [10] E. U. H. Qazi, T. Zia, and A. Almorjan, "Deep learning-based digital image forgery detection system," *Applied Sciences*, vol. 12, no. 6, 2022.
- [11] P. Ranjan, S. Patil, and F. Kazi, "Improved generalizability of deep-fakes detection using transfer learning based cnn framework," in *2020 3rd International Conference on Information and Computer Technologies (ICICT)*, pp. 86–90, 2020.
- [12] A. Diel, T. Lalg, F. S. Mellis, A. Teufel, and A. Bäuerle, "The harm of deepfakes: a scoping review of deepfakes' negative effects on human mind and behavior," *AI & Society*, Dec 2025.
- [13] A. Abdel-Wahab and M. Alkhatib, "Toward robust deepfake defense: A review of deepfake detection and prevention techniques in images," *Computers, Materials and Continua*, vol. 86, no. 2, pp. 1–34, 2025.
- [14] S. Zobaed, M. F. Rabby, M. Hossain, E. Hossain, M. S. Hasan, A. Karim, and K. Hasib, *DeepFakes: Detecting Forged and Synthetic Media Content Using Machine Learning*, pp. 177–201. 09 2021.
- [15] D. Epstein, I. Jain, O. Wang, and R. Zhang, "Online detection of ai-generated images," pp. 382–392, 10 2023.
- [16] Y. Patel, S. Tanwar, P. Bhattacharya, R. Gupta, T. Alsuwian, I. Davidson, and T. Mazibuko, "An improved dense cnn architecture for deepfake image detection," *IEEE Access*, vol. PP, 03 2023.