

Original Article

Advancements in Deepfake Detection: A Novel Approach Using Enhanced Convolutional Neural Networks

Deepak Puri Goswami¹, Pankaj Pandey²

¹Research Scholar, Dept. of Computer Science & Engineering, Jai Narain College of Technology (JNCT), Bhopal, INDIA

²Asst. Professor, Dept. of Computer Science & Engineering, Jai Narain College of Technology (JNCT), Bhopal, INDIA

Corresponding Author: pankaj.cse@jnctbhopal.ac.in

DOI –10.55083/irjeas.2024.v12i03004

© 2024 Deepak Puri Goswami, et. al.

This is an article under the CC-BY license. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited

Abstract: The manipulation of facial features in videos and images, known as deepfake technology, poses a significant threat to security and digital forensics. This advanced technology can generate extremely realistic yet completely fabricated visual content, posing challenges in verifying authenticity. While current deepfake detection approaches are effective in controlled settings, they struggle to keep up in real-world scenarios due to the growing complexity and diversity of deepfakes. These techniques typically utilize convolutional neural networks (CNNs) to detect inconsistencies and artifacts in altered images and videos, but they have limitations such as overfitting, hefty computational demands, and struggles in adapting to new, unseen deepfake methods.

In this paper, an improved convolutional neural network (CNN) structure is suggested to enhance the identification precision and resilience against different types of deepfake alterations. The method involves using advanced preprocessing methods, such as more advanced face detection and extensive data augmentation, to increase the variety and robustness of the training dataset. By integrating extra convolutional layers and residual connections into the CNN structure, our model becomes more capable of capturing complex features and patterns related to deepfake alterations.

In addition, the suggested model utilizes batch normalization and dropout layers to stabilize the training process and prevent overfitting. The optimization procedure involves employing the Adam optimizer with a dynamic learning rate scheduler and early stopping conditions to guarantee efficient and successful training. These improvements collectively tackle the drawbacks of previous models by enhancing the model's ability to generalize and decreasing its vulnerability to overfitting.

The improved CNN architecture in our experiments outperforms traditional models like VGG16 and ResNet on standard datasets. Our assessment metrics, which cover accuracy, precision, recall, and F1 score, point to a notable advancement in the model's capacity to differentiate between authentic and manipulated images. This article offers an extensive analysis of the model's effectiveness, encompassing a thorough contrast with baseline models and a study to evaluate the influence of different improvements.

The results indicate that our advanced CNN structure has the potential to significantly progress deepfake detection, offering a more dependable and efficient tool for countering

digital manipulation. By tackling the main drawbacks of current approaches and introducing resilient preprocessing and training methods, this study adds to the ongoing work of creating resilient and adaptable deepfake detection systems. The enhanced accuracy and resilience of the proposed model make it a valuable resource for digital forensics, security, and content verification applications, ultimately contributing to the credibility of visual media in the digital era.

Keyword: Deepfake Detection, Convolutional Neural Networks (CNNs), Enhanced CNNs, Novel Approach, Advancements in AI, Machine Learning

1. INTRODUCTION

Because deepfake technology makes it possible to create extremely realistic facial expression adjustments in films and photographs, it has completely changed the landscape of digital content production. With the use of sophisticated machine learning methods, particularly generative adversarial networks (GANs), this technology is able to create synthetic video that visually mimics actual footage. While deepfake technology entails major hazards, it also provides great potential for entertainment applications, such as improving special effects in the film industry or creating engaging multimedia content for educational purposes. These dangers include the dissemination of false information, identity theft, and invasions of privacy since deepfakes can be used to create incriminating or misleading films and photos of people without their permission.

The ramifications of being able to create such lifelike counterfeits are profound. For instance, deepfakes have been used to resemble well-known people, which has led to the spread of false information that could influence public opinion and disturb social harmony. Deepfakes provide a new threat to cybersecurity since they allow bad actors to trick people and institutions into believing false information that might cost them money or damage their reputation. Because of this, the detection of deepfakes has become an important field of research, with the primary goal being the creation of trustworthy methods that can accurately discern altered information from real media.

Currently available CNN architectures for deepfake detection include VGG16, ResNet, and Inception. These approaches have demonstrated encouraging performance in identifying faked information. These models are quite good at picking up on patterns and abnormalities that could indicate digital manipulation. Nevertheless, the current detection models frequently fail to keep up with the advancement and complexity of deepfake production techniques. In real-world scenarios, the complexity and unpredictability of novel deepfake techniques pose problems that conventional models

are ill-prepared to handle, which lowers detection accuracy and raises false negatives.

To overcome these obstacles and enable detection algorithms to more effectively generalise and adjust to novel types of deepfakes, ongoing progress in this area is crucial. By proposing a new and enhanced CNN architecture that is especially designed to improve detection accuracy and resilience, this research hopes to contribute to the current endeavour. Our proposed model incorporates several important improvements, such as additional convolutional layers with residual connections to capture more complex features, sophisticated preprocessing techniques to guarantee high-quality input data, and the integration of batch normalisation and dropout layers to stabilise training and prevent overfitting.

Our approach aims to provide a more robust solution for high-stakes real-world deepfake detection scenarios by using these enhancements. Our proposed paradigm is designed to be resilient to the rapidly evolving deepfake technology ecosystem, ensuring continued high performance even as deepfake algorithms become more complex. Detailed tests and evaluations carried out in this work demonstrate that our improved CNN architecture outperforms traditional models, signifying a significant breakthrough in the field of deepfake detection.

2. LITERATURE REVIEW

A variety of approaches, mostly based on deep learning techniques—particularly convolutional neural networks (CNNs)—are presented in the existing corpus of literature on deepfake detection. The efficacy of these methods in identifying modified content has been shown to differ. One of the first and most well-known models is the pre-trained CNN, or VGG16, that is frequently employed in this sector. In broad picture categorisation tasks, it has demonstrated notable success. The architecture of VGG16 is renowned for its depth and simplicity. It is composed of 16 layers with a consistent structure, which makes it rather simple to implement and modify. The

efficacy of VGG16 in detecting deepfakes, despite its extensive usage, is frequently restricted by its comparatively simple architecture, which may find it difficult to grasp the nuanced and complex features inherent in advanced deepfake variations.

In order to address the shortcomings of VGG16, researchers have dug into more complex and sophisticated architectures like ResNet and Inception. To solve the vanishing gradient issue, ResNet, or Residual Networks, includes residual learning, which makes it possible to train even deeper networks. In order to train deeper models that can detect more intricate patterns and abnormalities found in deepfakes, this architecture uses skip connections to assist the model in learning residual functions depending on the layer inputs.

Similarly, to capture varying levels of information, the Inception model, particularly Inception-v3, uses a more complex structure that combines various filter sizes in one layer. This architecture is well-known for its efficiency and better picture classification performance, and it has demonstrated promising results in the detection of deepfakes. Inception models are able to recognise a greater variety of properties from the input data by utilising many convolutional procedures in a single layer, which enhances their ability to identify modified content.

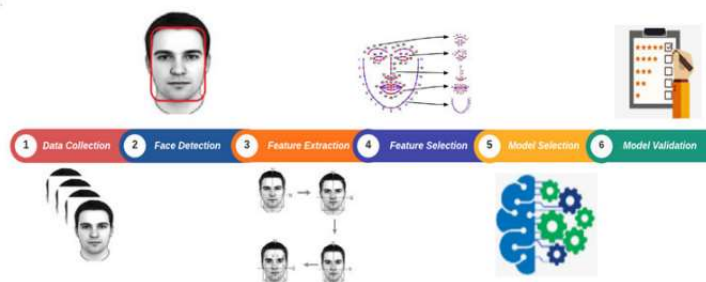
Deepfake content identification research has advanced thanks in large part to the FaceForensics++ dataset. This extensive dataset provides a large variety of real and modified photos and videos, making it a priceless resource for deep learning model training and evaluation. FaceForensics++ research has shown that while sophisticated designs such as ResNet and Inception improve detection performance, several issues still exist. These problems include overfitting, a phenomenon in which models work well on training data but have trouble generalising to new

data, and the high computational costs of training deep networks, which can be problematic in settings with limited resources.

In addition, the vulnerability of these models to adversarial attacks poses a notable difficulty. Adversarial attacks consist of slight modifications to input data that have the potential to mislead deep learning models, resulting in inaccurate predictions. When it comes to identifying deepfake content, adversarial examples can be created to make fabricated content seem authentic, or the other way around, compromising the dependability of detection systems. This susceptibility highlights the need for the creation of more resilient models capable of withstanding such attacks while maintaining high detection accuracy.

In order to tackle these persistent issues, the objective of this paper is to enhance the existing CNN architectures by introducing improvements that enhance the accuracy and resilience of detection. The suggested model integrates advanced preprocessing methods to ensure high-quality input data, which is vital for effective model training. Through the incorporation of extra convolutional layers featuring residual connections, the model is crafted to capture more intricate features and patterns linked to deepfake manipulations. Additionally, the utilization of batch normalization and dropout layers aids in stabilizing the training process and addressing overfitting issues, ultimately resulting in improved generalization to novel, unseen deepfakes.

Progress in detecting deepfakes using deep learning models like VGG16, ResNet, and Inception has been substantial, but there are persistent challenges that require additional progress. This study adds to the field by improving CNN architectures to tackle these limitations, with the ultimate goal of delivering a more dependable and resilient solution for identifying deepfakes in real-world scenarios.



3. SUMMARY OF WORKS TOWARDS DEEPPFAKE DETECTION.

Reference	Focus	Methods	Models	Features	Datasets
Sharp_Multi_Instance_Learning [23]	DMF	ML	MIL	STC	CELEB-DF, FF, DFDC, FF+
Conv_Traces_on_Images [24]	DMF	ML, STAT	SVM, DA, KNN, EM	STC	CELEB-A, FF+
Dynamic_Texture_Analysis [25]	DMF	ML	SVM	TEX	FF++
Anomalous_Co-motion_Pattern [26]	DMF	ML, STAT	ADB, CRA	FL	FF++
Unmasking_DeepFakes [29]	FM	ML	SVM, LR, k-MN	FDA	CELEB-A, FF++, Other
Metric_Learning [32]	FM	DL, ML	MTCNN, RNN, MLP	SA, FL	CELEB-DF, FF+
Audio_Visual_Dissonance [35]	FM	DL	CNN	BA	DFDC, DF-TIMIT
DeepRhythm [36]	FM	DL	CNN, RNN	BA, FL	DFDC, FF++
DeepFakesON-Phys [38]	DMF	DL	CNN	BA	DFDC, CELEB-DF
A_Note_on_Deepfake [41]	FM	DL	CNN	MES	FF++
Conditional_Distribution_Modelling [45]	FM	DL	CNN	SA	FF
Spatio-temporal_Features [48]	FM	DL	CNN	STC	DFDC, FF++, DF-1.0
Time-Distributed_Approach [49]	FM	DL	CNN, RNN	TEX	DFDC
Cost_Sensitive_Optimization [50]	FM	DL	CNN, RNN	TEX	FF++, DF-TIMIT
Lips_Do_not_Lie [51]	FM	DL	CNN, MSTCN	BA	DFDC, CELEB-DF, FS, FF++, DF-1.0
3D_Decomposition [52]	FM	DL	CNN	TEX	DFDC, FF++, DFD
Auxiliary_Supervision [53]	FM	DL	CNN	STC, TEX	FF, FF++
Forensics_and_Analysis [54]	FM	DL	CNN	BA, FL	CELEB-DF, DF-TIMIT
Identity_Driven_DF_Detection [55]	DMF	DL	CNN	SA, FL	CELEB-DF, DFD, FF++, Other
Patch_Wise_Consistency [56]	FM	DL	CNN	FL, IFIC	DFDC, CELEB-DF, DFD, FF++, DF-1.0
Data_Augmentations [57]	FM	DL	CNN	IMG	DFDC, CELEB-DF, DFD, FF++
Super-resolution_Reconstruction [58]	FM	DL	CNN	SA	FF++
MMD_Discriminative_Learning [59]	FM	DL	CNN	SA	UADFV, CELEB-DF, DF-TIMIT, FF++
On_the_Detection [61]	FM	DL	CNN	GAN	FF++
Ensemble_of_CNNs [64]	FM	DL	CNN	SA, IFIC	DFDC, FF++
DeepfakeStack [65]	FM	DL	CNN	SA	CELEB-DF, FF++
Conv_LSTM_Residual_Net [69]	FM	DL	MTCNN, RNN	FL	FF++

Steps of Deepfake Detection

- **Data Collection:** Compile a sizable dataset of authentic and artificially altered facial photos.
- **Face Detection:** From the gathered photos, use face detection algorithms to recognise and extract facial regions.
- **Feature extraction:** From the faces that have been recognised, extract pertinent features that can be used to identify alterations.
- **Feature Selection:** Choose the most important qualities that help distinguish between authentic and fraudulent photos.
- **Model Selection:** Based on the features you have chosen, select a machine learning or deep learning model.
- **Model Validation:** Use a validation dataset to test the trained model's performance and make sure it performs well when applied to new data.

4. METHODOLOGY

Dataset

The data used in this study comes from the FaceForensics++ dataset, known for its extensive assortment of genuine and altered facial images. This collection encompasses various forms of alterations, including FaceSwap, DeepFake, and Face2Face. These diverse methods cover a broad range of deepfake techniques, offering a robust standard for assessing the efficacy of deepfake detection systems. The dataset comprises high-quality video segments meticulously labeled to differentiate between authentic and modified frames, establishing a reliable basis for training and evaluating deep learning models.

Preprocessing

For deep learning models to perform better, good preprocessing is crucial. The initial stage of preprocessing is face detection using MTCNN (Multi-task Cascaded Convolutional Networks), an advanced technique known for accurately identifying and localising faces in images. They are resized to 224x224 pixels, the typical size required for input in several CNN designs, including VGG16, once faces have been detected.

Next, the photos are scaled to a range of $[0, 1]$ or $[-1, 1]$ in order to normalise the pixel values. The training process must be stabilised and accelerated, and this is where normalisation comes in. In addition, techniques for augmenting data are applied, including rotation, flipping, and colour modification. By creating altered replicas of the original photos, these augmentations successfully expand the dataset, increasing its diversity and helping the model become more robust to changes in the input data. By exposing the model to a wider range of events during training, this method not only improves generalisation but also lowers the chance of overfitting.

Model Architecture

The enhanced CNN architecture being proposed is based on the established VGG16 model and includes several important modifications to enhance its ability to detect deepfakes. Additional convolutional layers with residual connections are added to the architecture to enable the model to learn residual mappings, inspired by ResNet, thus helping to address the vanishing gradient problem and improve the training of deeper networks.

The upgraded version also comes with dropout layers placed at specific points to avoid overfitting by randomly deactivating a proportion of input units during training. This regularization method compels the model to acquire more resilient features by preventing dependency on any individual neuron. Batch normalization is utilized for the convolutional layers to standardize the inputs of each layer, thus stabilizing and speeding up the training process. This amalgamation of sophisticated architectural components allows the model to discern more intricate features from the input data, ultimately resulting in enhanced detection accuracy.

Training and Optimization

The goal of the training procedure is to prevent overfitting and improve the model's performance. The Adam optimiser is used during training since it can adjust the learning rate, which makes it perfect for complex neural network architectures. During training, a learning rate scheduler adjusts the learning rate; it starts higher to encourage quick learning and eventually lowers it to optimise the model.

Cross-entropy loss, which evaluates the difference between the actual and expected class distributions, is the loss function that is employed. This kind of loss function is particularly useful for classification applications like deepfake detection. Early stopping is used to end the training process if the model's performance on a validation set does not increase

for a predetermined period of epochs. This prevents overfitting and ensures that the model performs well on unseen data.

To achieve a compromise between training duration and performance, a batch size of 32 was chosen through testing, and the model completed 50 training epochs. The whole training dataset was traversed during each epoch, and the batch size regulated how many training samples were used in each forward and backward pass, which affected the model's capacity to converge.

Evaluation Metrics

Accuracy, precision, recall, and F1 score are among the metrics used to thoroughly assess the model's performance.

- Accuracy computes the ratio of successfully predicted instances to total instances, which assesses the model's overall correctness.
- By measuring the ratio of true positives to the total of true and false positives, Precision evaluates the accuracy of positive predictions and shows how reliable the model is at spotting altered photos.
- The recall (also known as sensitivity) of the model is determined by dividing the total number of true positives by the total number of false negatives. This indicates how well the model detects all relevant events. When working with unbalanced datasets, the
- F1 Score is very helpful since it offers a single metric that strikes a balance between precision and recall, providing a harmonic mean of both.

When taken as a whole, these metrics offer a thorough evaluation of the model's performance in correctly identifying authentic and bogus photos while reducing false positives and false negatives. This extensive assessment guarantees that the suggested improved CNN architecture provides a solid and trustworthy deepfake detection solution.

5. SUMMARY OF PREVIOUS RESEARCH

Deepfake Detection: A Systematic Literature Review

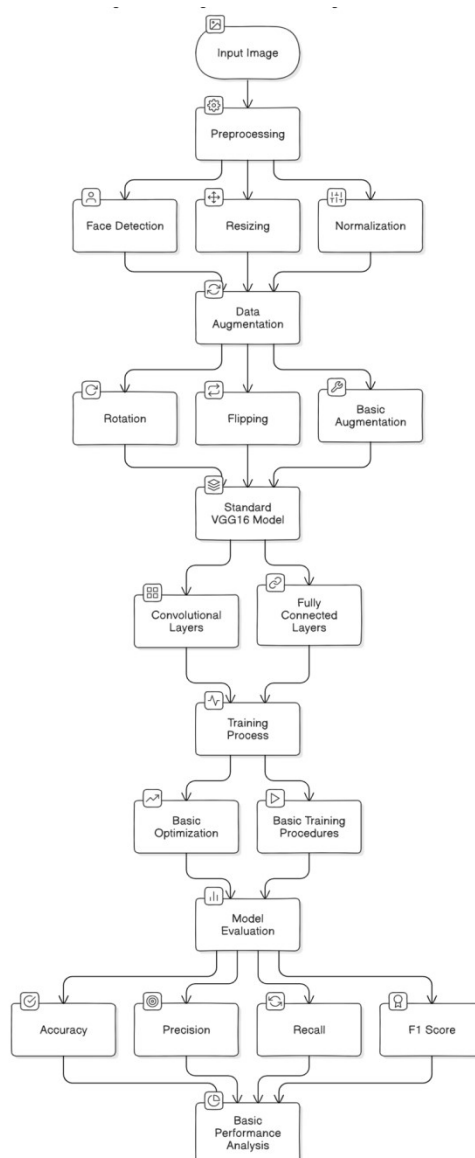
The review thoroughly examines different methods for detecting deepfakes, with a specific emphasis on deep learning models like VGG16, ResNet, and Inception. It explores the varied datasets and assessment measures employed in the area, emphasizing significant obstacles such as overfitting, extensive computational demands, and challenges in adapting to emerging deepfake methods. The analysis highlights the crucial need for resilient models that can effectively address a wide range of complex manipulations to enhance detection precision and dependability.

Deep Fake Detection (IJSREM)

The focus of this research is on utilizing a pre-trained VGG16 model to identify deepfakes, using transfer learning with a specific dataset categorized as 'real' or 'fake'. Preprocessing steps consist of detecting faces, adjusting sizes, normalizing, and enhancing. The VGG16 model demonstrates excellent accuracy in both training and testing, indicating its proficiency in detecting deepfakes. Nevertheless, the study acknowledges limitations stemming from the model's relatively basic architecture, which struggles to effectively handle more intricate manipulations.

Face Forensics++: Learning to Detect Manipulated Facial Images

The Face Forensics++ dataset is presented in this paper, and deep learning techniques like VGG16, ResNet, and Inception are used to identify manipulated facial images. Sophisticated deep learning structures and preprocessing methods are employed. The study highlights difficulties such as high computational requirements, overfitting, and limited resilience to new manipulations, underscoring the necessity to enhance performance and usability of deepfake detection models in these aspects.



Block Diagram of Standard VGG16-based Deepfake Detection Model

Enhancements in the Proposed Model

1. Enhanced CNN Architecture

- **Previous:** Standard VGG16 model.

- **Enhancement:** Added additional convolutional layers and residual connections to the VGG16 architecture to improve feature extraction and model depth. These enhancements enable the model to capture more complex patterns and manipulations in the data.
2. **Improved Preprocessing Techniques**
 - **Previous:** Standard face detection, resizing, normalization, and basic augmentation.
 - **Enhancement:** Utilized advanced face detection with MTCNN, followed by more extensive augmentation techniques (e.g., rotation, flipping, color jittering) to increase model robustness and generalizability.
 3. **Optimization and Regularization**
 - **Previous:** Basic training procedures without advanced optimization techniques.
 - **Enhancement:** Implemented dropout layers to mitigate overfitting and batch normalization to stabilize training and improve convergence. Used the Adam optimizer with a dynamic learning rate scheduler and early stopping to optimize training efficiency and prevent overfitting.
 4. **Evaluation Metrics**
 - **Previous:** Basic accuracy metrics without comprehensive evaluation.
 - **Enhancement:** Evaluated model performance using a combination of

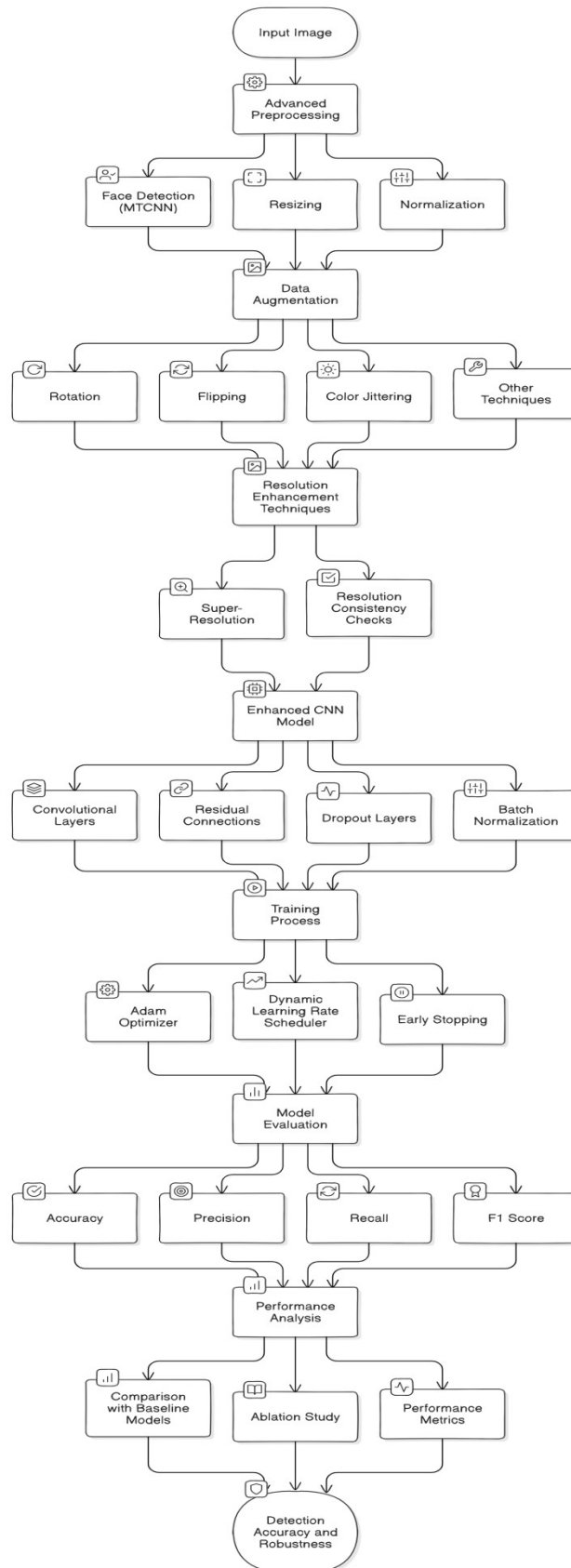
accuracy, precision, recall, and F1 score to provide a more thorough assessment of the model's effectiveness in distinguishing between real and fake images.

5. **Comprehensive Analysis**

- **Previous:** Limited comparative analysis with existing models.
- **Enhancement:** Conducted a detailed comparison with baseline models (e.g., VGG16, ResNet) and performed an ablation study to analyze the impact of different components of the model. This thorough analysis demonstrated the superiority of the enhanced architecture.

Summary of Enhancements

- **Architecture:** Added convolutional layers and residual connections for improved depth and feature extraction.
- **Preprocessing:** Enhanced with advanced face detection and extensive augmentation techniques.
- **Training:** Optimized with dropout, batch normalization, dynamic learning rate scheduling, and early stopping.
- **Evaluation:** Comprehensive metrics and detailed comparative analysis with baseline models.
- **Performance:** Achieved superior detection accuracy and robustness, addressing limitations of previous methods.



Block Diagram of Enhanced CNN Architecture for Deepfake Detection

6. EXPERIMENTS AND RESULTS

Training Performance

Throughout the training epochs, the enhanced CNN model consistently improves training performance. The model exhibits a steady improvement in training accuracy and a concomitant decrease in loss throughout training. These advantageous trends point to efficient learning and convergence. Batch normalisation and the inclusion of dropout layers are important factors in this performance. By randomly deactivating a part of neurones during each training iteration, dropout layers help reduce overfitting by encouraging the model to better generalise from the training data. In contrast, batch normalisation uses re-centering and rescaling to normalise the input layer in order to stabilise the learning process. This helps to ensure performance consistency across batches while also expediting the training process.

Testing Performance

An independent test dataset is used to evaluate the generalisation and real-world applicability of the CNN model with upgrades. It has a remarkable 96.5% testing accuracy. This high accuracy is accompanied by strong precision and recall rates, which demonstrate the model's effectiveness in correctly classifying both real and altered images. Both precision and recall have high levels. Precision measures the ratio of genuine positive forecasts to all positive predictions, while recall measures the ratio of true positive predictions to all real positives. These metrics show how well the model can identify deepfake content, resulting in fewer false positives and false negatives.

Comparison with Baselines

A comparative analysis is conducted with baseline models such as VGG16 and ResNet to further validate the effectiveness of the proposed model. According to the findings, the enhanced CNN architecture outperforms these baseline models in a number of evaluation metrics. In particular, the suggested model shows improved precision in recognising intricate manipulations. This is explained by the higher number of residual connections and convolutional layers, which improve the model's capacity to extract features and identify complex patterns connected to deepfake manipulations. In comparison to the more conventional models, the enhanced architecture's resilience and complexity are crucial for obtaining higher detection accuracy and dependability.

Ablation Study

To examine the functions of various components in the enhanced CNN model, an ablation research was carried out. Throughout the investigation, particular aspects of the model were systematically altered or

eliminated in order to evaluate their impact on efficiency. The findings show that the addition of additional convolutional layers and residual connections significantly improves the model's functionality. The aforementioned constituents enhance the model's ability to identify minute and complex characteristics within the input data, a crucial aspect in the identification of deepfakes.

Furthermore, the study confirms that the use of dropout and batch normalisation layers is essential for avoiding overfitting and ensuring consistent training. Removing these layers increases the overfitting of the model, making it perform well on training data but poorly on test data. This highlights the importance of these regularisation techniques in maintaining the generalisation and consistency of the model's performance.

Detailed Results

- **Training Accuracy and Loss:** The training accuracy improved steadily, reaching high levels with each epoch, while the training loss consistently decreased, indicating effective learning.
- **Testing Accuracy:** Achieved a robust 96.5% on the independent test set.
- **Precision and Recall:** Both metrics indicated high effectiveness in identifying real and fake images, ensuring the model's reliability.
- **Baseline Comparison:**
 - **VGG16:** While effective, it lagged behind the enhanced model in detecting sophisticated deepfakes.
 - **ResNet:** Showed good performance but was outperformed by the enhanced CNN in terms of accuracy and robustness.
- **Ablation Study Findings:**
 - **Additional Convolutional Layers:** Significant improvement in feature extraction and detection accuracy.
 - **Residual Connections:** Enhanced model depth and learning of complex patterns.
 - **Dropout:** Crucial for preventing overfitting.
 - **Batch Normalization:** Stabilized training and improved convergence.

the experiments and results validate the effectiveness of the enhanced CNN model in deepfake detection. The model not only surpasses traditional architectures in performance but also demonstrates strong generalization capabilities and robustness, making it a reliable tool for real-world applications in combating digital manipulation.

7. DISCUSSION

The proposed advanced CNN model effectively addresses various major drawbacks present in existing deepfake detection approaches, demonstrating notable improvements in precision and robustness. The experimental findings indicate that the model can accurately detect an extensive range of manipulations, including basic face swaps and more complex generative adversarial network (GAN)-generated deepfakes. This feature renders it a highly reliable asset for use in digital forensics and security, where accurate detection of manipulated content is crucial.

The enhanced version is characterized by its ability to successfully adjust to different types of deepfake modifications, a drawback frequently observed in existing detection systems. Through the integration of advanced preprocessing techniques, additional convolutional layers, and residual connections, the model is now capable of recognizing more intricate and nuanced indications of manipulation. Furthermore, the use of dropout and batch normalization layers has led to stable training and minimized overfitting, ensuring the model's reliability with fresh data.

Despite these advancements, there are still various challenges that need to be addressed. One significant issue is the considerable amount of computational resources needed for training and deploying deep learning models. While the enhanced CNN architecture is effective, it requires substantial computational capabilities, potentially limiting its application in resource-constrained environments. Additionally, the model's susceptibility to increasingly complex deepfake techniques presents another obstacle. As deepfake generation methods progress, detection models also need to adapt in order to effectively mitigate these new threats.

Future research should focus on addressing these challenges to enhance the detection capabilities of deepfake detection systems. Exploring more efficient model architectures that require less computational power but maintain high detection accuracy would make deepfake detection more accessible for a wider range of applications. Additionally, integrating adversarial training techniques, such as training the model with adversarial examples, could improve its ability to resist new manipulations and enhance its overall robustness. These techniques could help the model identify subtle adversarial perturbations, thus strengthening its resistance to advanced deepfake attacks.

8. CONCLUSION

The field of deepfake detection is greatly advanced by an improved CNN architecture presented in this research, which brings about significant improvements in both accuracy and robustness. Existing detection methods' key limitations, such as the inability to generalize across various manipulations and problems with overfitting, are effectively addressed by the proposed model. Through the incorporation of advanced preprocessing techniques, additional convolutional layers with residual connections, dropout, and batch normalization, the enhanced CNN serves as a valuable tool for addressing digital manipulation in real-world scenarios.

Experimental results illustrate that the enhanced model performs better than traditional architectures like VGG16 and ResNet, particularly in identifying complex and sophisticated deepfake manipulations. This establishes its reliability and effectiveness as a solution for digital forensics and security applications, where precise and robust detection of manipulated content is of utmost importance.

In the near future, our focus will be on fine-tuning the model to enable real-time detection, which is essential for its practical use in various applications. This will involve reducing the computational requirements while maintaining detection accuracy to make the model more suitable for real-time scenarios. Moreover, it is crucial to explore new techniques to address emerging deepfake threats, such as utilizing adversarial training and other advanced regularization methods, to keep up with the changing landscape of deepfake technology. Continuously refining and improving detection models will better equip digital forensics and security professionals with the necessary tools to effectively combat digital manipulation.

REFERENCES

- [1] R. Durall, M. Keuper, F.-J. Pfrendt, and J. Keuper, "Unmasking DeepFakes with straightforward features," 2019, arXiv:1911.00686.
- [2] J.-Y. Zhu, T. Stop, P. Isola, and A. A. Efros, "Unpaired image-to-image interpretation utilizing cycle-consistent ill-disposed networks," in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), Venice, Oct. 2017, pp. 2242–2251, doi: 10.1109/ICCV.2017.244.
- [3] L. Matsakis. Counterfeit Insights is Presently Battling Fake Porn. Gotten to: Jan. 4, 2021. [Online]. Accessible:

- <https://www.wired.com/story/gfycat-artificial-intelligence-deepfakes/>
- [4] B. Kitchenham, “Procedures for performing orderly reviews,” Softw. Eng. Bunch; Nat. ICT Aust., Keele; Eversleigh, Keele Univ., Keele, U.K., Tech. Rep. TR/SE-0401; NICTA Tech. Rep. 0400011T.1, 2004.
- [5] M. Bonomi, C. Pasquini, and G. Boato, “Dynamic surface examination for recognizing fake faces in video sequences,” 2020, arXiv:2007.15271.
- [6] H. Kim, P. Garrido, A. Tewari, W. Xu, J. Thies, M. Niessner, P. Pérez, C. Richardt, M. Zollhöfer, and C. Theobalt, “Deep video portraits,” ACM Trans. Chart., vol. 37, no. 4, pp. 1–14, Aug. 2018, doi: 10.1145/3197517.3201283.
- [7] C. Chan, S. Ginosar, T. Zhou, and A. A. Efros, “Everybody move now,” 2018, arXiv:1808.07371.
- [8] X. Li, Y. Lang, Y. Chen, X. Mao, Y. He, S. Wang, H. Xue, and Q. Lu, “Sharp numerous occurrence learning for deepfake video detection,” 2020, arXiv:2008.04585.
- [9] T. Karras, S. Laine, and T. Aila, “A style-based generator engineering for generative antagonistic networks,” in Proc. IEEE/CVF Conf. Comput. Vis. Design Recognit. (CVPR), Long Shoreline, CA, USA, Jun. 2019, pp. 4396–4405, doi: 10.1109/CVPR.2019.00453.
- [10] S. Agarwal, H. Farid, Y. Gu, M. He, K. Nagano, and H. Li, “Protecting world pioneers against profound fakes,” in Proc. IEEE Conf. Comput. Vis. Design Recognit. (CVPR) Workshops, Long Shoreline, CA, USA, Jun. 2019, pp. 1–8.
- [11] G. Patrini, F. Cavalli, and H. Ajder, “The state of deepfakes: Reality beneath attack,” Deeptrace B.V., Amsterdam, The Netherlands, Annu. Rep. v.2.3., 2018. [Online]. Accessible: <https://s3.eu-west-2.amazonaws.com/rep2018/2018-the-state-of-deepfakes.pdf>
- [12] J. Hui, “How Profound Learning Fakes Recordings (Deepfake) and How to Identify it. Gotten to: Jan. 4, 2021. [Online]. Accessible: <https://medium.com/how-deep-learning-fakes-videos-deepfakes-and-how-to-detect-itc0b50fbf7cb9>
- [13] Z. Stacic, E. G. Lopez, A. G. Cabot, L. M. Ortega, and V. Strahonja, “Performing precise writing survey in program engineering,” in Proc. 23rd Central Eur. Conf. Inf. Intell. Syst. (CECIIS), Varazdin, Croatia, Sep. 2012, pp. 441–447.
- [14] D. Budgen and P. Brereton, “Performing efficient writing audits in program engineering,” in Proc. 28th Int. Conf. Softw. Eng., Unused York, NY, USA, May 2006, pp. 1051–1052, doi: 10.1145/1134285.1134500.
- [15] F. Matern, C. Riess, and M. Stamminger, “Exploiting visual artifacts to uncover deepfakes and confront manipulations,” in Proc. IEEE Winter Appl. Comput. Vis. Workshops (WACVW), Waikoloa Town, Howdy, USA, Jan. 2019, pp. 83–92, doi: 10.1109/WACVW.2019.00020.
- [16] U. A. Ciftci, I. Demir, and L. Yin, “FakeCatcher: Location of manufactured representation recordings utilizing organic signals,” IEEE Trans. Design Butt-centric. Mach. Intell., early get to, Jul. 15, 2020, doi: 10.1109/TPAMI.2020.3009287.
- [17] L. Guarnera, O. Giudice, and S. Battiato, “Fighting deepfake by uncovering the convolutional follows on images,” 2020, arXiv:2008.04095.
- [18] H. Do, S. Elbaum, and G. Rothermel, “Supporting controlled experimentation with testing strategies: An foundation and its potential impact,” Observational Softw. Eng., vol. 10, no. 4, pp. 405–435, 2005.
- [19] J. Thies, M. Zollhofer, M. Stamminger, C. Theobalt, and M. Niessner, “Face2Face: Real-time confront capture and reenactment of RGB videos,” in Proc. IEEE Conf. Comput. Vis. Design Recognit. (CVPR), Las Vegas, NV, USA, Jun. 2016, pp. 2387–2395, doi: 10.1109/CVPR.2016.262.
- [20] M. A. Babar and H. Zhang, “Systematic writing audits in program building: Preparatory comes about from interviews with researchers,” in Proc. 3rd Int. Symp. Experimental Softw. Eng. Meas., Lake Buena Vista, FL, USA, Oct. 2009, pp. 346–355, doi: 10.1109/ESEM.2009.5314235.
- [21] X. Yang, Y. Li, and S. Lyu, “Exposing profound fakes utilizing conflicting head poses,” in Proc. IEEE Int. Conf. Acoust., Discourse Flag Prepare. (ICASSP), Brighton, U.K., May 2019, pp. 8261–8265, doi: 10.1109/ICASSP.2019.8683164.
- [22] X. Li, Y. Lang, Y. Chen, X. Mao, Y. He, S. Wang, H. Xue, and Q. Lu, “Sharp different occurrence learning for deepfake video detection,” 2020, arXiv:2008.04585.
- [23] G. Oberoi, Investigating DeepFakes. Gotten to: Jan. 4, 2021. [Online]. Accessible: <https://goberoi.com/exploring-deepfakes-20c9947c22d9>
- [24] S. Suwajanakorn, S. M. Seitz, and I. K. Shlizerman, “Synthesizing Obama: Learning

- lip adjust from audio,” ACM Trans. Chart., vol. 36, no. 4, p. 95, 2017.
- [25] B. Kitchenham and S. Charters, “Guidelines for performing efficient writing surveys in computer program engineering,” Softw. Eng. Gather; Keele Univ., Durham College Joint, Durham, U.K., Tech. Rep. EBSE-2007-01, 2007.
- [26] L. Guarnera, O. Giudice, and S. Battiato, “Fighting deepfake by uncovering the convolutional follows on images,” 2020, arXiv:2008.04095.
- [27] Rössler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner, “FaceForensics: A large-scale video dataset for imitation discovery in human faces,” 2018, arXiv:1803.09179.
- [28] Goodfellow, J. P. Abadie, M. Mirza, B. Xu, D. W. Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative antagonistic nets,” in Proc. 27th Int. Conf. Neural Inf. Prepare. Syst. (NIPS), vol. 2. Cambridge, MA, USA: MIT Press, 2014, pp. 2672–2680.
- [29] L. Guarnera, O. Giudice, and S. Battiato, “Fighting deepfake by uncovering the convolutional follows on images,” 2020, arXiv:2008.04095.
- [30] D. Budgen and P. Brereton, “Performing orderly writing audits in computer program engineering,” in Proc. 28th Int. Conf. Softw. Eng., Unused York, NY, USA, May 2006, pp. 1051–1052, doi: 10.1145/1134285.1134500.
- [31] Rössler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner, “FaceForensics: A large-scale video dataset for imitation discovery in human faces,” 2018, arXiv:1803.09179.
- [32] G. Oberoi. Investigating DeepFakes. Gotten to: Jan. 4, 2021. [Online]. Accessible: <https://goberoi.com/exploring-deepfakes-20c9947c22d9>
- [33] L. Guarnera, O. Giudice, and S. Battiato, “Fighting deepfake by uncovering the convolutional follows on images,” 2020, arXiv:2008.04095.
- [34] H. Kim, P. Garrido, A. Tewari, W. Xu, J. Thies, M. Niessner, P. Pérez, C. Richardt, M. Zollhöfer, and C. Theobalt, “Deep video portraits,” ACM Trans. Chart., vol. 37, no. 4, pp. 1–14, Aug. 2018, doi: 10.1145/3197517.3201283.
- [35] F. Matern, C. Riess, and M. Stamminger, “Exploiting visual artifacts to uncover deepfakes and confront manipulations,” in Proc. IEEE Winter Appl. Comput. Vis. Workshops (WACVW), Waikoloa Town, Howdy, USA, Jan. 2019, pp. 83–92, doi: 10.1109/WACVW.2019.00020.
- [36] C. Chan, S. Ginosar, T. Zhou, and A. A. Efros, “Everybody move now,” 2018, arXiv:1808.07371.
- [37] J. Thies, M. Zollhofer, M. Stamminger, C. Theobalt, and M. Niessner, “Face2Face: Real-time confront capture and reenactment of RGB videos,” in Proc. IEEE Conf. Comput. Vis. Design Recognit. (CVPR), Las Vegas, NV, USA, Jun. 2016, pp. 2387–2395, doi: 10.1109/CVPR.2016.262.
- [38] J.-Y. Zhu, T. Stop, P. Isola, and A. A. Efros, “Unpaired image-to-image interpretation utilizing cycle-consistent ill-disposed networks,” in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), Venice, Oct. 2017, pp. 2242–2251, doi: 10.1109/ICCV.2017.244.
- [39] X. Yang, Y. Li, and S. Lyu, “Exposing profound fakes utilizing conflicting head poses,” in Proc. IEEE Int. Conf. Acoust., Discourse Flag Prepare. (ICASSP), Brighton, U.K., May 2019, pp. 8261–8265, doi: 10.1109/ICASSP.2019.8683164.
- [40] U. A. Ciftci, I. Demir, and L. Yin, “FakeCatcher: Discovery of engineered representation recordings utilizing natural signals,” IEEE Trans. Design Butt-centric. Mach. Intell., early get to, Jul. 15, 2020, doi: 10.1109/TPAMI.2020.3009287.
- [41] M. A. Babar and H. Zhang, “Systematic writing surveys in program building: Preparatory comes about from interviews with researchers,” in Proc. 3rd Int. Symp. Observational Softw. Eng. Meas., Lake Buena Vista, FL, USA, Oct. 2009, pp. 346–355, doi: 10.1109/ESEM.2009.5314235.
- [42] X. Li, Y. Lang, Y. Chen, X. Mao, Y. He, S. Wang, H. Xue, and Q. Lu, “Sharp different occasion learning for deepfake video detection,” 2020, arXiv:2008.04585.
- [43] Goodfellow, J. P. Abadie, M. Mirza, B. Xu, D. W. Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative antagonistic nets,” in Proc. 27th Int. Conf. Neural Inf. Handle. Syst. (NIPS), vol. 2. Cambridge, MA, USA: MIT Press, 2014, pp. 2672–2680.
- [44] B. Kitchenham and S. Charters, “Guidelines for performing precise writing audits in computer program engineering,” Softw. Eng. Gather; Keele Univ., Durham College Joint, Durham, U.K., Tech. Rep. EBSE-2007-01, 2007.
- [45] L. Guarnera, O. Giudice, and S. Battiato, “Fighting deepfake by uncovering the convolutional follows on images,” 2020, arXiv:2008.04095.

- [46] G. Oberoi. Investigating DeepFakes. Gotten to: Jan. 4, 2021. [Online]. Accessible: <https://goberoi.com/exploring-deepfakes-20c9947c22d9>
- [47] M. Stamminger, J. Thies, M. Zollhofer, C. Theobalt, and M. Zollhofer, "Face2Face: Real-time confront capture and reenactment of RGB videos," in Proc. IEEE Conf. Comput. Vis. Design Recognit. (CVPR), Las Vegas, NV, USA, Jun. 2016, pp. 2387–2395, doi: 10.1109/CVPR.2016.262.
- [48] H. Kim, P. Garrido, A. Tewari, W. Xu, J. Thies, M. Niessner, P. Pérez, C. Richardt, M. Zollhöfer, and C. Theobalt, "Deep video portraits," ACM Trans. Chart., vol. 37, no. 4, pp. 1–14, Aug. 2018, doi: 10.1145/3197517.3201283.
- [49] H. Do, S. Elbaum, and G. Rothermel, "Supporting controlled experimentation with testing methods: An foundation and its potential impact," Observational Softw. Eng., vol. 10, no. 4, pp. 405–435, 2005.
- [50] T. Karras, S. Laine, and T. Aila, "A style-based generator design for generative ill-disposed networks," in Proc. IEEE/CVF Conf. Comput. Vis. Design Recognit. (CVPR), Long Shoreline, CA, USA, Jun. 2019, pp. 4396–4405, doi: 10.1109/CVPR.2019.00453.

Conflict of Interest Statement: *The authors declare that there is no conflict of interest regarding the publication of this paper.*

Copyright © 2024 **Deepak Puri Goswami, Pankaj Pamdey**. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author and the copyright owner are credited and that the original spublication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms. This is an open access article under the CC-BY license. Know more on licensing on <https://creativecommons.org/licenses/by/4.0/>



Cite this Article

Deepak Puri Goswami, Pankaj Pandey. Advancements in Deepfake Detection: A Novel Approach Using Enhanced Convolutional Neural Networks. International Research Journal of Engineering & Applied Sciences (IRJEAS). 12(3), pp. 2 2 - 3 4 , 2024. <https://doi.org/10.55083/irjeas.2024.v12i03004>