# Progress Report

July 2024

**Title:** *Topic to be decided*

*Progress Report submitted in partial fulfillment of the requirement for the award of the degree.*
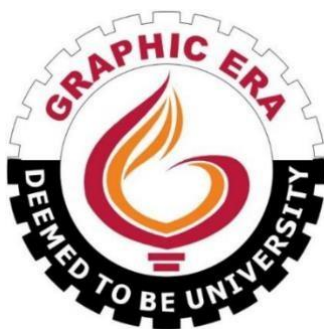
of

## Doctor of Philosophy

in

## Computer Science and Engineering

By

## YOGESH LOHUMI
**(EN NO. PHD23CSE139)**

Under the supervision of
## DR. D. R. GANGODKAR



# DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

Graphic Era (Deemed to Be) University, Dehradun – 248002

July, 2024

## CANDIDATE'S DECLARATION

I hereby certify that the work carried out for the thesis entitled …… *Topic to be decided (Coursework Completed*) …… in partial fulfillment of the requirements for the award of the degree of Doctor of Philosophy is a bonafide original research work by me under the supervision of **Dr. D. R. Gangodkar** during the semester (Jan 2024 – June 2024) at department of **Computer Science and Engineering** in lab on more than 100 working days. This work has not been proposed elsewhere for the award of a degree/diploma/certificate.

**Signature of the Candidate**

This is to certify that the above-mentioned statement in the candidate's declaration is correct to the best of my/our knowledge.

Date: _____

**Signature of Supervisor**

# ABSTRACT

Deepfake video has usefulness in entertainment and multimedia technology, however, the danger of deepfake is significant to the social, economic, and political sectors so far. Specifically, to diverge any public opinion by generating fake news and spreading misleading information, national security may be under risk due to misrepresenting statements given by political leaders. The creation of such manipulated videos are getting easier day by day and at the same time it is necessary to detect and prevent them. Deepfakes refer to synthetic media where a person in an existing image or video is replaced with someone else's likeness. The challenge of detecting deepfakes is crucial due to their potential for misuse. Deepfakes refer to media—such as videos, images, and audio—where an individual's likeness is replaced with someone else's using artificial intelligence (AI) techniques. The rise of deepfakes has led to significant concerns regarding misinformation, privacy, and cybersecurity. Detecting deepfakes is crucial for preventing their malicious use.

In order to do that, researchers are creating challenging fake video databases for artificial intelligence (AI) based detection models to contribute to the research. Deepfakes are a recent off-the-shelf manipulation technique that allows anyone to swap two identities in a single video. In addition to Deepfakes, a variety of GAN based face swapping methods have also been published with accompanying code. To counter this emerging threat, we have constructed an extremely large face swap video dataset to enable the training of detection models, and organized the accompanying DeepFake Detection Challenge (DFDC) Kaggle competition. Importantly, all recorded subjects agreed to participate in and have their likenesses modified during the construction of the faceswapped dataset. The DFDC dataset is by far the largest currentlyand publicly-available face swap video dataset, with over 100,000 total clips sourced from 3,426 paid actors, produced with several Deepfake, GAN-based, and non-learned methods. In addition to describing the methods used to construct the dataset, we provide a detailed analysis of the top submissions from the Kaggle contest. We show although Deepfake detection is extremely difficult and still an unsolved problem, a Deepfake detection model trained only on the DFDC can generalize to real "in-the-wild" Deepfake videos, and such a model can be a valuable analysis tool when analyzing potentially Deepfaked videos. With this understanding, we did not construct our dataset from publicly available videos. Instead, we commissioned a set of videos to be taken of individuals who agreed to be filmed, to appear in a machine learning dataset, and to have their face images manipulated by machine learning models. In order to reflect the potential harm of Deepfaked videos designed to harm a single, possibly non-public person, videos were shot in a variety of natural settings without professional lighting or makeup, (but with high-resolution cameras, as resolution can be easily downgraded).

**Keywords:** Machine Learning, Deep learning, DeepFake, CNNs, GANs, Encoder

# TABLE OF CONTENTS

# ✟ WORK DONE

In the first semester, the academic journey was shaped by carefully selected coursework that provided both theoretical knowledge and practical skills. The courses included:

- Research Methodology (PHDM – 102): Provided an organized framework for methodological investigation.
- Machine Learning (MCS – 128): Examined models and algorithms central to computational intelligence.
- Applied Data Science (MCS – 131): Introduced basic theory and practical applications of Data Science.
- Seminar on RM Presentation (SEM – 001): Focused on intellectual exploration, communication, and research presentation skills.

The examinations for these courses were conducted in February 2024 – March 2024, with the following results:

- Research Methodology (PHDM – 102): 72 (Grade S)
- Machine Learning (MCS – 128): 78 (Grade S)
- Applied Data Science (MCS – 131): 85 (Grade S)
- Seminar on RM Presentation (SEM – 001): 42 (Grade S)

# ✟ WORK IN PROGRESS

The first semester has seen the successful completion of all required coursework, including examinations, which has provided a solid theoretical and practical foundation in key areas.

Currently, the focus is on identifying and refining specific research objectives that will guide the upcoming phases of the research. Concurrently, active efforts are being made in the preparation of a detailed research proposal, outlining the objectives, methodologies, and expected outcomes. Additionally, significant progress has been made in conducting a comprehensive literature review, which highlights the current state of Deepfake Detection research, key trends, applications, and identifies significant gaps and opportunities for further exploration.

This report primarily details the work done on the literature review, setting the stage for the focused research ahead.

# CHAPTER 1

# <u>INTRODUCTION</u>

Deep Learning (DL) has been effectively utilized in various complicated challenges in healthcare, industry, and academia for various purposes, including thy- roid diagnosis, lung nodule recognition, computer vision, large data analytics, and human-level control. Nevertheless, developments in digital technology have been used to produce software that poses a threat to democracy, national security, and confidentiality. Deepfake is one of those DL-powered apps that has lately sur- faced. So, deepfake systems can create fake images primarily by replacement of scenes or images, movies, and sounds that humans cannot tell apart from real ones. Various technologies have brought the capacity to change a synthetic speech, image, or video to our fingers. Furthermore, video and image frauds are now so convincing that it is hard to distinguish between false and authentic con- tent with the naked eye. It might result in various issues and ranging from deceiv- ing public opinion to using doctored evidence in a court. For such considerations, it is critical to have technologies that can assist us in discerning reality. This study gives a complete assessment of the literature on deepfake detection strategies using DL-based algorithms. We categorize deepfake detection methods in this work based on their applications, which include video detection, image detection, audio detection, and hybrid multimedia detection. Considering the time continuity in videos, Guera et al. [1] first proposed to use RNN to detect deepfake videos. In their work, autoencoders was found to be completely unaware of previously generated faces because faces were generated frame-by-frame. This lack of temporal awareness results in multiple anomalies, whichare crucial evidence for deepfake detection. To check the continuity between adjacent frames, an end-to-end trainable recurrent deepfake video detection system was proposed. As Figure 6 shows, the proposed system is mainly composed of a convolutional long short- term memory (LSTM) structure for processing frame sequences. Two essential components are used in a con- volutional LSTM structure, where CNN is used for frame feature extraction and LSTM is used for temporal sequence analysis. Specifically, a pretrained inceptionV3 [2] is adapted to output a deep representation for each frame. The 2048- dimensional feature vectors extracted by the last pooling layers are applied as the sequential LSTM input, character- izing the continuity between image sequences. Finally, a fullyconnected layer and a softmax layer are added to compute forgery probabilities of the frame sequence tested. The experiments on a self-made dataset showed that the al- gorithm can accurately make predictions even when the length of a video is less than 2 s. Although this research did not show its superiority since there were no large-scale datasets at the time, several articles after were inspired. After the time-based detection method showed its

effective- ness, many related studies were proposed. In [3], Sabir et al. utilized the temporal information present in the video stream to detect deepfake videos. Similar to [1], an endto-end model is built, where CNN is also involved in the follow-up training. Meanwhile, face alignment based on facial landmarks and spatial transformer network is applied to further improve the performance of the algorithm. Even though such solutions guarantee high accuracy in videos with high quality, they do not perform well on low-quality video when the continuity be- tween adjacent frames is disrupted by video compression op- erations. To solve this problem, a CNN-RNN framework based on automatic weighting mechanisms was proposed by Montserrat et al. [4]. Considering that the face qualities of some frames are not high, an automatic weighting mechanism was proposed to emphasize the most reliable regions when making a video-level prediction. Experiments showed that combining CNN and RNN achieves high detection accuracies on the DFDC dataset. Except for the robustness of algo- rithms, generalization ability is also essential for forgery detection tasks. Zhao et al. [5] used optical flow to capture the obvious differences of facial expressions between adjacent frames. However, these studies did not show strong general- ization or robustness. To solve this problem, Wu et al. [6] proposed a novel manipulation detection framework, named SSTNet, exploiting both low-level artefacts and temporal dis- crepancies. Another study proposed by Masi et al. [7] ob- tained good generalization on multiple datasets. In their research, a two-branch recurrent network is applied to prop- agate the original information while suppresses the face con- tent. Multiband frequencies are amplified using a Laplacian of Gaussian as a bottleneck layer. Inspired by [8], a new loss function is designed for better isolating manipulated face. The experimental results on several datasets show the excellent generalization performance of the detection algorithm.

# CHAPTER 2

## <u>TIMELINE CHART</u>

The timeline chart outlines the key milestones and progression from Coursework Examination to the of initial exploration of the topics. Each box in Figure 1 represents a significant step in the timeline.
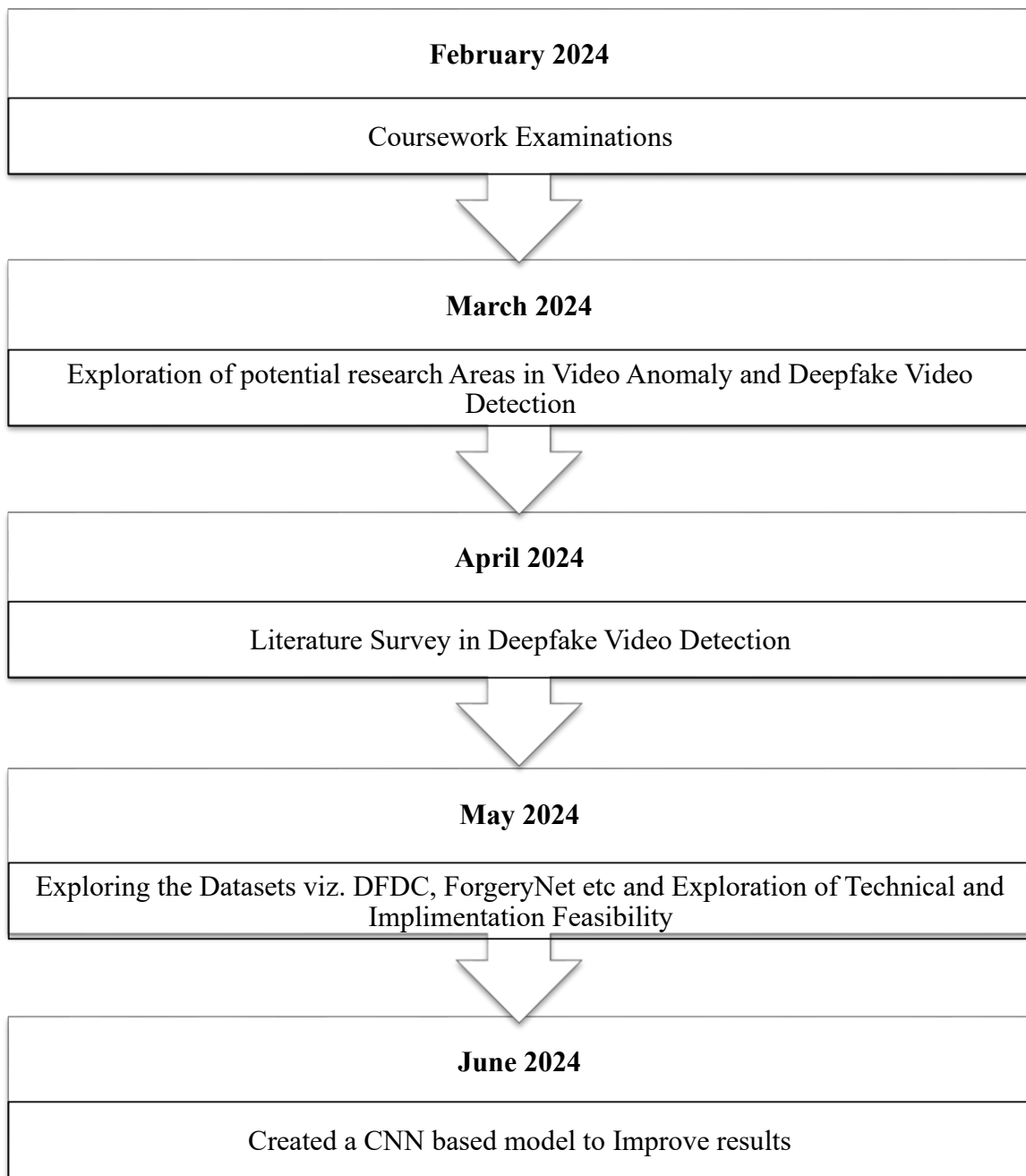


| **February 2024** |
| :---: |
| Coursework Examinations |

| **March 2024** |
| :---: |
| Exploration of potential research Areas in Video Anomaly and Deepfake Video Detection |

| **April 2024** |
| :---: |
| Literature Survey in Deepfake Video Detection |

| **May 2024** |
| :---: |
| Exploring the Datasets viz. DFDC, ForgeryNet etc and Exploration of Technical and Implimentation Feasibility |

| **June 2024** |
| :---: |
| Created a CNN based model to Improve results |

Fig. 1 PhD Second Semester Timeline.

# CHAPTER 3

## <u>DETAILS OF THE WORK DONE</u>

This chapter offers an in-depth examination of the endeavors undertaken during the designated timeframe, encompassing coursework and initial forays into research. The section begins by elucidating the coursework subjects selected and the insights gained from each. It then delves into a brief literature review conducted in DeepFake Detection, highlighting key findings and trends.

3.1 LITERATURE REVIEW

3.1 <u>Deepfake Detection</u>

In recent years, video forensic techniques have received more attention due to the development of video forgery techniques. Researchers have proposed many forensic methods for video forgery, such as moving objects detection [18], removing objects detection [19], copy-move detection [20], frame interpolation detection [21], frame deletion detection [22] etc. However, the above methods fail to detect deepfake videos because deepfake videos are generated by artificial intelligence facial manipulation technology. The initial manipulated facial video has many apparent defects. Some detection methods are based on handcrafted features, highlighting specific failures during the generation process. For example, some methods rely on eye blinking [8] and facial expression change [7], which has a specific frequency and duration in humans and cannot be replicated in deepfake videos. Li and Lyu [23] found that the image blending operation would leave face warp artifacts and used CNN models to extract these artifacts. In [9], the authors note that some forgery methods have a common step: blending the forged face with the existing background, and they performed the detection by predicting the blending boundary of the fake face. Yang et al. [24] observed that the forgery process of splicing the synthesized face region into the original image would introduce head poses errors, and they used the difference in estimated head poses as a feature vector to detect deepfake. The advanced Deepfake technology overcomes these defects and becomes more visually realistic. Some works [10], [11] has extracted features from the spatial domain for classification and achieved excellent performance on specific datasets. Afchar et al. [25] presented Meso-4 and MesoInception-4 networks with few layers to capture the mesoscopic properties of images. In [13], the authors proposed a method that leverages the capsule network to detect face manipulation. In [26], Zhao et al. proposed a

multiattentional deepfake detection network, and they believe that the difference between real and fake images is often subtle and local. To improve the robustness of the detection methods to postprocessing operations, Chen et al. [27] designed a dual-stream network by integrating dual-color spaces RGB and YCbCr and using an improved Xception model. Yu et al. [15] proposed a commonality learning framework for detecting similar traces left by different forgery methods to improve generalization when detecting unknown forgery methods. Recently, some researchers have considered the temporal information of videos. In [14], a convolutional Long Short Term Memory (LSTM) network is used to exploit such dependencies and improve upon singleframe analysis. In [16], Zhang et al. adapted 3DCNN to capture inconsistent features between frames. Hu et al. [28]. adopted a two streams network for compressed deepfake videos detection, the frame-level stream with a low complexity network to extract spatial features, and the temporalitylevel stream to extract the inconsistency between frames. Chen et al. proposed Xception-LSTM network using a spatiotemporal attention mechanism, improved performance by enhancing spatiotemporal correlation and frame structure information. The complex deep convolutional neural networks efficiently handle classification problems in the pixel domain. In addition, some methods focus on the frequency domain to capture artifacts of fake faces. For example, Durall et al. [29] used Discrete Fourier Transform (DFT) to extract frequency-domain information. F3-Net [5] took advantage of frequency-aware decomposed image components and analyzed the statistic features. In [30], Liu et al. proposed a spatial-phase shallow learning (SPSL) method that combines spatial image and phase spectrum to capture the up-sampling artifacts of face forgery. Chen et al. [31] proposed a multi-scale patch similarity module, fusing information in both RGB and frequency domains for local feature representation. Most existing methods treat deepfake detection as a general binary classification problem. However, they tend to suffer from overfitting. To solve this problem, we propose novel artifacts disentangle framework to improve detection performance by removing irrelevant information. B. Disentanglement Learning Disentanglement learning has achieved great success in several tasks in computer vision. In face recognition tasks, feature representation of multiple attributes becomes a challenge, [32] disentangles face into identity and pose vectors for pose-invariant face recognition. TDGAN [33] learns to disentangle expressional information from other unrelated facial attributes for facial expression recognition. In [34], not only the 2D facial images but also the 3D face reconstruction disentangles the representation of a 3D face into identity, expressions, poses, albedo, and illuminations. In gait recognition, Zhang et al. [35] disentangles the representations of appearance, canonical, and pose features from an input gait video. Video prediction also uses the

idea of disentangled representation. Denton et al. [36] proposed DRNet, where representations are disentangled into content and pose, and the poses are penalized for encoding semantic information with the use of a discrimination loss. Similarly, MCNet [37] disentangled motion from content using image differences and shared a single content vector in prediction. In face anti-spoofing, [38] disentangled the spoof traces from input faces and implemented generic detection of various spoof types, such as print, replay, mask, mannequin head, etc. In this paper, different from, the artifacts in deepfake are much weaker and are concentrated in smaller areas in some forgery methods. Artifact extraction in deepfake is more susceptible to interference from face identity and background information, and the proposed ADAL aims to eliminate the influence of irrelevant features.

Deepfake detection methods fall into three categories [34, 37]. Methods in the first category focus on the physical or psychological behavior of the videos, such as tracking eye blinking or head pose movement. The second category focus on GAN fingerprint and biological signals found in images, such as blood flow that can be detected in an image. The third category focus on visual artifacts. Methods that focus on visual artifacts are data-driven, and require a large amount of data for training, the implemented falls into the third category. In this section, we will discuss various architectures designed and developed to detect visual artifacts of Deepfakes. Darius et al. [1] proposed a CNN model called MesoNet network to automatically detect hyper-realistic forged videos created using Deepfake [40] and Face2Face [54]. The authors used two network architectures (Meso-4 and MesoInception-4) that focus on the mesoscopic properties of an image. Yuezun and Siwei [34] proposed a CNN architecture that takes advantage of the image transform (i.e., scaling, rotation and shearing) inconsistencies created during the creation of Deepfakes. Their approach targets the artifacts in affine face warping as the distinctive feature to distinguish real and fake images. Their method compares the Deepfake face region with that of the neighboring pixels to spot resolution inconsistencies that occur during face warping. Huy et al. [41] proposed a novel deep learning approach to detect forged images and videos. The authors focused on replay attacks, face swapping, facial reenactments and fully computer-generated image spoofing. Daniel Mas Montserrat et al. [38] proposed a system that extracts visual and temporal features from faces present in a video. Their method combines a CNN and RNN architecture to detect Deepfake videos. Md. Shohel Rana and Andrew H. Sung [50] proposed a Deepfake Stack, an ensemble method (A stack of different DL models) for Deepfake detection.

The ensemble is composed of XceptionNet, InceptionV3, InceptionResNetV2, MobileNet, ResNet101, DenseNet121, and DenseNet169 open-source DL models. Junyaup Kim et al. [29] proposed a classifier that distinguishes target individuals from a set of similar people using ShallowNet, VGG-16, and Xception pre-trained DL models. The main objective of their system is to evaluate the classification performance of the three DL models. 3. Convolutional Vision Transformer In this section, we present our approach to detect Deepfake videos. The Deepfake video detection model consists of two components: the preprocessing component and the detection component. The preprocessing component consists of the face extraction and data augmentation. The detection components consist of the training component, the validation component, and the testing component. The training and validation components contain a Convolutional Vision Transformer (CViT). The CViT has a feature learning component that learns the features of input images and a ViT architecture that determines whether a specific video is fake or real. The testing component applies the CViT learning model on input images to detect Deepfakes. Our proposed model is shown in Figure 2
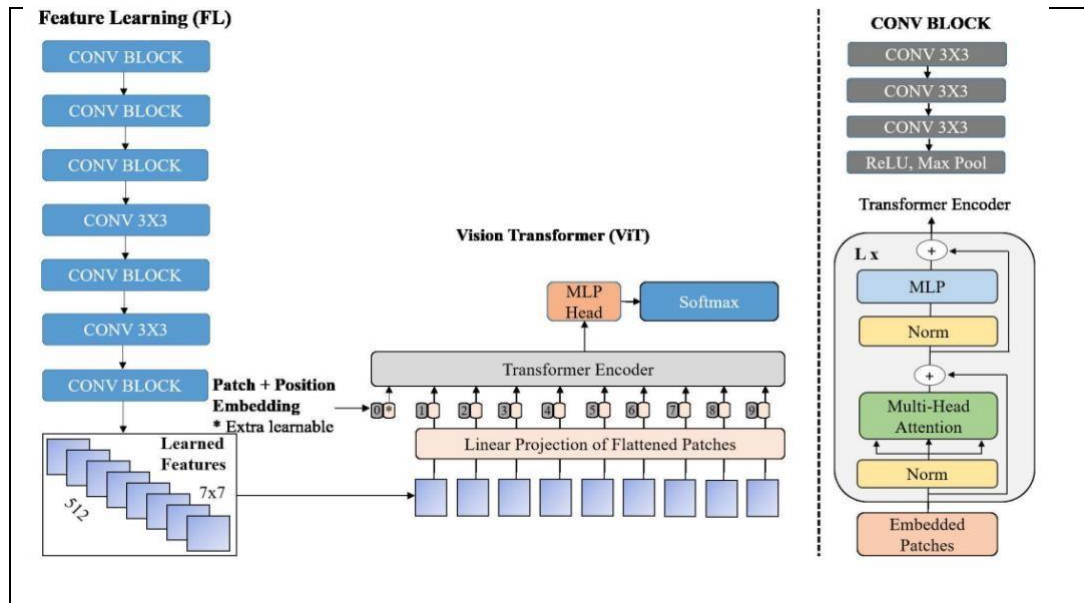


Fig.2 Implemented CNN model

3.2 Explanation of Model Layers

1.  Convolutional Layers: These layers apply convolution operations to the input data, extracting features such as edges, textures, and patterns.

2.  Max-Pooling Layers: These layers reduce the spatial dimensions of the feature maps, retaining the most important features and reducing computational complexity.

3.  Batch Normalization: This layer normalizes the activations of the previous layer, improving training speed and stability.

4.  Flatten Layer: This layer converts the 2D feature maps into a 1D vector, preparing the data for the fully connected layers.

5.  Dense Layers: Fully connected layers that learn complex representations by combining features from previous layers.

6.  Dropout Layers: These layers randomly drop a fraction of the neurons during training, preventing overfitting.

Compiling the Model

After defining the model architecture, the next step is to compile it. Compiling the model involves specifying the optimizer, loss function, and metrics to use during training.

# CHAPTER 4

## <u>RESULTS AND CONCLUSION</u>

### 4.1    Training Accuracy Comparison

We compared the performance of two models: a custom CNN model and a fine-tuned ResNet50 model, for detecting deepfakes using a dataset of face images as shown in Table 1

Table. 1 Training Accuracy Comparision

| Model | Training Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|
| Custom CNN | 100.00% | 0.6594 | 0.6070 | 0.6321 |
| Fine-tuned ResNet50 | 59.76% | 0.5842 | 0.5290 | 0.5552 |

### 4.2 Performance Metrics
- Custom CNN Model:

    o Training Accuracy: The custom CNN achieved a perfect training accuracy of 100.00%, indicating it learned to classify the training data perfectly. o Precision: It correctly identified 65.94% of predicted deepfakes. o Recall: The model captured 60.70% of actual deepfakes in the dataset.

    o F1-score: Achieved an F1-score of 0.6321, balancing precision and recall.

- Fine-tuned ResNet50 Model:

    o Training Accuracy: The fine-tuned ResNet50 achieved a training accuracy of 59.76%, indicating it correctly classified 59.76% of the training data.

    o Precision: It correctly identified 58.42% of predicted deepfakes.

    o Recall: The model captured 52.90% of actual deepfakes in the dataset.

    o F1-score: Achieved an F1-score of 0.5552, showing a balanced performance between precision and recall.

Based on the training accuracy comparison and performance metrics, the custom CNN model shows promise for detecting deepfakes with high accuracy on the training data. Further optimization and fine-tuning of parameters could potentially enhance its performance. The fine-tuned ResNet50 model, despite its lower training accuracy, demonstrates reasonable performance metrics, suggesting

its capability to identify deepfakes to a certain extent. Fine-tuning hyperparameters and possibly increasing the training data size could improve its effectiveness.

In conclusion, this comprehensive approach covers the entire process of building, training, evaluating, and optionally fine-tuning a deep learning model for detecting deepfakes. Each step from data preprocessing and model compilation to training and evaluation is crucial for developing a robust and accurate deepfake detection system. By leveraging advanced techniques such as data augmentation and fine-tuning pre-trained models, we can enhance the model's performance and ensure its effectiveness in real-world applications. This project underscores the importance of leveraging AI and machine learning to address emerging challenges in cybersecurity and digital forensics.

## 4.3 PLAN FOR NEXT SEMESTER

The emphasis for the following semester will be on developing a research proposal, establishing the framework for an insightful and guided exploration of various research areas in Computer Vision

References

1. D. Güera and E. J. Delp, "Deepfake Video Detection Using Recurrent Neural Networks," *2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, Auckland, New Zealand, 2018, pp. 1-6

2. C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens and Z. Wojna, "Rethinking the Inception Architecture for Computer Vision," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, 2016, pp. 2818-2826

3. E. Sabir, J. Cheng, A. Jaiswal, W. AbdAlmageed, I. Masi, and P. Natarajan, "Recurrent convolutional strategies for face manipulation detection in videos," arXiv [cs.CV], 2019.

4. D. Mas Montserrat, H. Hao, S. K. Yarlagadda, S. Baireddy, R. Shao, J. Horváth, et al., "Deepfakes detection with automatic face weighting", arXiv:2004.12027, 2020.

5. S. Ha, M. Kersner, B. Kim, S. Seo, and D. Kim, "MarioNETte: Few-shot face reenactment preserving identity of unseen targets," Proc. Conf. AAAI Artif. Intell., vol. 34, no. 07, pp. 10893–10900, 2020.

6. Y. Lu, J. Chai, and X. Cao, "Live speech portraits: Real-time photorealistic talking-head animation," arXiv [cs.GR], 2021.

7. X. Wu, Z. Xie, Y. Gao and Y. Xiao, "SSTNet: Detecting Manipulated Faces Through Spatial, Steganalysis and Temporal Features," *In Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Barcelona, Spain, 2020, pp. 2952-2956.

8. X. Wu, Z. Xie, Y. Gao and Y. Xiao, "SSTNet: Detecting Manipulated Faces Through Spatial, Steganalysis and Temporal Features," In Proc. *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Barcelona, Spain, 2020, pp. 2952-2956,

9. H. Farid, "Image forgery detection," IEEE Signal Process. Mag., vol. 26, no. 2, pp. 16–25, Mar. 2009.

10. I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in Proc. Adv. Neural Inf. Process. Syst., vol. 27, 2014, pp. 1–9.

11. P. Baldi, "Autoencoders, unsupervised learning, and deep architectures," in Proc. ICML Workshop Unsupervised Transf. Learn., 2012, pp. 37–49.

12. T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2019, pp. 4401–4410.

13. Y. Mirsky and W. Lee, "The creation and detection of deepfakes: A survey," ACM Comput. Surv., vol. 54, no. 1, pp. 1–41, Jan. 2022.

14. M. Masood, M. Nawaz, K. M. Malik, A. Javed, and A. Irtaza, "Deepfakes generation and detection: State-of-the-art, open challenges, countermeasures, and way forward," 2021, arXiv:2103.00484.

15. R. Tolosana, R. Vera-Rodriguez, J. Fierrez, A. Morales, and J. Ortega-Garcia, "Deepfakes and beyond: A survey of face manipulation and fake detection," Inf. Fusion, vol. 64, pp. 131–148, Dec. 2020.

16. T. T. Nguyen, Q. V. H. Nguyen, D. T. Nguyen, D. T. Nguyen, T. Huynh-The, S. Nahavandi, T. T. Nguyen, Q.-V. Pham, and C. M. Nguyen, "Deep learning for deepfakes creation and detection: A survey," 2019, arXiv:1909.11573.

17. L. Verdoliva, "Media forensics and DeepFakes: An overview," IEEE J. Sel. Topics Signal Process., vol. 14, no. 5, pp. 910–932, Aug. 2020.

18. K. Fukushima, "Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position," Biol. Cybern., vol. 36, no. 4, pp. 193–202, Apr. 1980.

19. Y. LeCun, B. Boser, J. Denker, D. Henderson, R. Howard, W. Hubbard, and L. Jackel, "Handwritten digit recognition with a back-propagation network," in Proc. Adv. Neural Inf. Process. Syst., vol. 2, 1989, pp. 396–404.

20. Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," Proc. IEEE, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.

21. J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu, X. Wang, G. Wang, J. Cai, and T. Chen, "Recent advances in convolutional neural networks," Pattern Recognit., vol. 77, pp. 354–377, May 2018.

22. S. Dong, P. Wang, and K. Abbas, "A survey on deep learning and its applications," Comput. Sci. Rev., vol. 40, May 2021, Art. no. 100379.

23. O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, and A. C. Berg, "ImageNet large scale visual recognition challenge," Int. J. Comput. Vis., vol. 115, no. 3, pp. 211–252, Dec. 2015.

24. M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in Proc. Eur. Conf. Comput. Vis., Cham, Switzerland: Springer, 2014, pp. 818–833.

25. K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014, arXiv:1409.1556.

26. C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2015, pp. 1–9.

27. K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2016, pp. 770–778.

28. S. Hochreiter and J. Schmidhuber, "Long short-term memory," Neural Comput., vol. 9, no. 8, pp. 1735–1780, 1997.

29. A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," 2015, arXiv:1511.06434.

30. A. Creswell and A. A. Bharath, "Inverting the generator of a generative adversarial network," IEEE Trans. Neural Netw. Learn. Syst., vol. 30, no. 7, pp. 1967–1974, Jul. 2019.

31. T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of GANs for improved quality, stability, and variation," 2017, arXiv:1710.10196.

32. A. Brock, J. Donahue, and K. Simonyan, "Large scale GAN training for high fidelity natural image synthesis," 2018, arXiv:1809.11096.

33. T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila, "Analyzing and improving the image quality of StyleGAN," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2020, pp. 8110–8119.

34. T. Karras, M. Aittala, J. Hellsten, S. Laine, J. Lehtinen, and T. Aila, "Training generative adversarial networks with limited data," 2020, arXiv:2006.06676.

35. Generated Photos. Face Generator—Generate Faces Online Using AI. [Online]. Available: https://generated.photos/face-generator

36. H. Dang, F. Liu, J. Stehouwer, X. Liu, and A. K. Jain, "On the detection of digital face manipulation," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2020, pp. 5781–5790.

37. J. C. Neves, R. Tolosana, R. Vera-Rodriguez, V. Lopes, H. Proenca, and J. Fierrez, "GANprintR: Improved fakes and evaluation of the state of the art in face manipulation detection," IEEE J. Sel. Topics Signal Process., vol. 14, no. 5, pp. 1038–1048, Aug. 2020.

38. Y. Zhu, Q. Li, J. Wang, C. Xu, and Z. Sun, "One shot face swapping on megapixels," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2021, pp. 4834–4844.

39. Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo, "StarGAN: Unified generative adversarial networks for multi-domain image-to-image translation," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., Jun. 2018, pp. 8789–8797.

40. Z. He, W. Zuo, M. Kan, S. Shan, and X. Chen, "AttGAN: Facial attribute editing by only changing what you want," IEEE Trans. Image Process., vol. 28, no. 11, pp. 5464–5478, Nov. 2019.