

# Deep Learning for the Generation and Detection of Deep Fakes

<sup>1</sup> Pradyumna Yambair, <sup>2</sup> MG. Suma

<sup>1</sup> Assistant Professor, Megha Institute of Engineering & Technology for Women, Ghatkesar.

<sup>2</sup> MCA Student, Megha Institute of Engineering & Technology for Women, Ghatkesar.

## Article Info

Received: 30-04-2025

Revised: 06 -06-2025

Accepted: 17-06-2025

Published:28/06/2025

## Abstract—

Computer vision, image identification, and natural language processing are just a few of the many fields that have made use of deep learning. The development of deep learning algorithms for picture recognition and modification has resulted in the emergence of deepfakes, which use these techniques to generate synthetic pictures that may sometimes be difficult to differentiate from the originals. Deepfake image detection methods have proliferated in response to growing worries about individuals' right to privacy and security. This paper investigates these methods and suggests using deep learning to improve the quality of deepfakes made.

Index Terms—deepfake, deep learning, Artificial intelligence, machine learning, tensorflow

## INTRODUCTION

Deepfake is an application that emerged from machine vision, which is making steady progress in several domains, including basic image identification software, robotics, and automotive [1] [2] [3]. The term "deepfake" refers to a method that use deep learning algorithms to generate seemingly authentic-looking fake photographs. This is achieved by, for example, replacing one person's face in a source image with another's in a target image. Using deep learning encoders and decoders, which have been widely employed in the machine vision area [4] [5], is the fundamental technique for creating deepfakes. The process begins with the en coders extracting all of the picture's characteristics, and then the false image is generated using decoders. Nowadays, it's easy to find large datasets of images on social media. This abundance of data has allowed for the development of more advanced deepfake techniques, as training the deep learning models used to be a challenging task. Deepfake algorithms are often built on top of Tensorflow [6]. TensorFlow is a free and open-source software framework for data-flow graph-based numerical computing. Although Google

created it for internal use in its own machine learning and deep neural network research and development, the system is broadly applicable and has found great popularity for machine learning applications since it was made publicly available and free to use. Because TensorFlow's APIs are compatible with Python, we can use them to build neural networks rapidly while still achieving sufficient performance. This allows us to experiment with different CNN architectures without touching a large amount of code. You could watch a video of a prominent public figure or the president giving a speech and not know if it's real or fake [7]. The process of creating these fake images and videos is much easier today; all you need is an image or video of the target individual to generate the fake content. This poses a serious threat to a future where fake news is everywhere. There are a rising amount of deepfakes on the internet, thus major tech firms are constantly studying strategies to identify them. A number of organizations have recently banded together to promote more investment in research and development for the detection and prevention of deepfakes. These include Facebook,

Amazon, Microsoft, and the Partnership on AI's Media Integrity Steering Committee [8]. In addition, Google has made a public dataset available for free as part of the deepfake detection challenge [9]. The fact that tech giants like Google and Microsoft are concerned about deepfake highlights the gravity of the problem. This research delves into one approach to identifying deepfake photos by using Mesonet CNN. Section II covers deepfake production, Section III covers deepfake detection, Section IV presents experimental results, and Section V offers concluding thoughts. This is the general outline of the study.

## Deep fake creation

By using deep learning techniques, the goal of creating a deepfake is to substitute the face of a certain individual in a video or picture with that of another person. Apps like FakeApp and FaceSwap, which are accessible online and made by developers and communities, are notable examples of such user-friendly programs. Some references: [10] [11]. In order to compress images, deepfake uses an autoencoder-decoder pipeline, which is common in the field of deep neural networks; by creating a bottleneck in the network, this pipeline compels the network to produce a compressed version of the original input [12]. High-quality picture compression is becoming attainable with the introduction of increasingly powerful encoders, which may make deepfake tasks easier by reducing the amount of CPU power needed [13] [14]. Two autoencoders are trained to create deepfakes. A deepfake image is created by first learning the features of the source image with an autoencoder and then learning the features of the target image with an encoder. The two encoders share their parameters and then use the decoder from the source image to reconstruct the target image, resulting in a picture of the target that has features from the source. Figure 1. According to DeepFaceLab [15] [16] and a plethora of other sources, this is the gold standard for creating deepfake images. Thanks to social media, there is an abundance of publicly available images and videos on platforms like Instagram and YouTube. This is particularly true for public figures and celebrities,

who were the initial targets of deepfakes and continue to be the most impacted by them [17]. The performance of this method is directly correlated to the size of the dataset available to train the deep neural network. There are a plethora of datasets available for usage in deepfake from the scientific community as well.

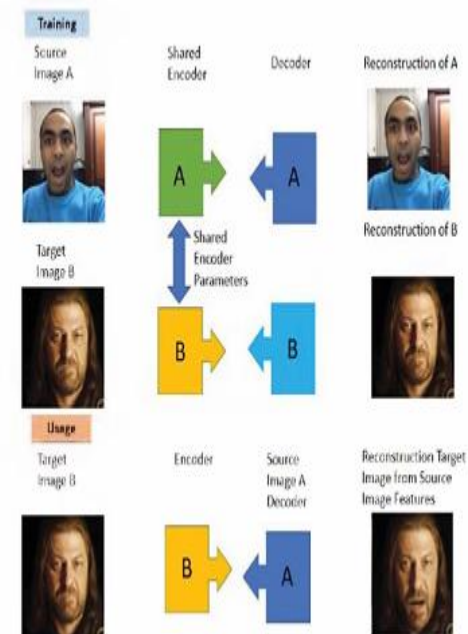


Fig. 1: Deepfake Creation Layout.

Using these concepts, one may create deepfake movies or photos; however, images are more efficient because of their tiny size and low processing requirements, making them quicker to manufacture. The proliferation of false news and the development of deepfake detection systems pose a serious threat. Figure 2 shows how user-friendly deepfake technologies are making it easy to create these photos, and how common they are as potential components of fake news stories. That is why, day by day, the importance of technologies that can identify deepfakes grows. Thirdly, the Deepfake Deposition One frequent usage of deepfake is to generate very difficult-to-detect swapped facial photos [18]. There are a plethora of ways for both making and finding deepfakes in the literature [19].



Fig. 2: Deepfake generated image.

Afchar et al. [20] suggested using MesoNet, a few-layer deep neural network. CNN showed an impressively high detection rate of over 98% for Deepfake. Part A. Convolutional Neural Networks A neural network developed with the express purpose of detecting deepfakes is Mesonet. Instagram videos, for example, are typically low-quality compressed videos, making microscopic analysis based on image noise impossible; MesoNet accounts for this; furthermore, detecting deepfakes at a higher semantic level is difficult; even humans have trouble doing it sometimes [21]. Consequently, MesoNet uses a small-layer deep neural network as an intermediary approach. Starting with a configuration of four layers of sequential convolutions and pooling (Fig. 3), this network then moves on to a dense network with a single hidden layer. It is usual practice to use a convolution layer first, then a pooling layer, when extracting image features; this is because the convolution layer finds the features and the pooling layer generates a down sampled version of the feature map.

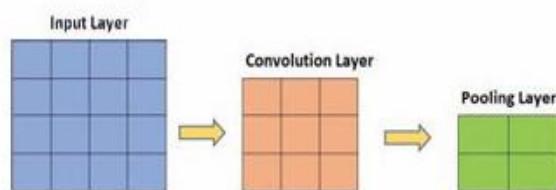


Fig. 3: Convolution and Pooling Layers

The convolutional layers use ReLU activation functions, which introduce non-linearities, and Batch Normalization [22] to regularize the output; the fully-connected layers, meanwhile, utilize Dropout to

regularize and enhance their resilience; and CNNs exploit this pattern to extract feature maps, which improves generalization. Section IV. Practical The outcomes A. The Making of Deepfakes While there are other Deepfake production methods documented in the literature, the two most prevalent include picture animation techniques and face swapping. First Order Motion Model for Image Animation is the name of the method that was used [23]. Figure 4 shows how the model accepts a source video and a target picture as inputs and uses the source video to animate the target image.



Fig. 4: deepfake video generated from source image

The Improving of Images using Deep Learning It was investigated if creating deepfakes using picture enhancement algorithms may increase their quality. The first stage was to test the single-image super resolution technique ESRGAN [24]. Although the output was a higher-resolution version of the input, the deepfake picture didn't improve significantly enough to warrant this additional step. The next step was to use DFDNet, a blind face restoration approach that yielded much superior results; the deepfakes that emerged from this process were noticeably higher in quality. Image 5. Using BasicSR, an open-source image and video restoration toolkit that includes features like super resolution, denoise, deblurring, JPEG artifacts removal, etc., these approaches were evaluated [26]. It offers a simple and fast method to evaluate several image enhancement techniques, like ESRGAN, DFDNet, StyleGAN2, and so on.



Fig. 5: Increasing quality of Deepfake with DFDNet

**Identifying Deep Fakes** A big collection of training photos is necessary to construct a strong neural network that can recognize sophisticated deepfakes. Companies like Google give big datasets to assist speed up the study on defense against deepfakes, and social media makes it easy to get these kinds of photos. More than 5,000 photos were utilized by MesoNet. The photos are categorized into two types: actual and deepfake (Fig. 6).



Fig. 6: dataset

Fig. 7 shows that after training the CNN with the dataset, it can recognize deepfake pictures with a confidence rating of above 80%. It is easier to spot deepfake pictures of a person's face gazing at an angle than one where the person's face is staring directly at the camera, however deepfake defects are readily apparent while perusing the existing deepfake photographs that are accessible online. When it comes to slanted faces, the existing deepfake creation methods aren't very good.



Fig. 7: Correctly classified images with high confidence rate

more than eighty-percent Nevertheless, even with a big dataset, some deepfake photos are still mistaken for genuine ones due to the invisibility of the images generated by sophisticated deepfake generation algorithms. Figure 8: Volume V: Concussion The use of deep learning for picture improvement and deepfake generation and detection were both investigated in this research.



means of enhancing the produced deepfakes. In contrast to the often poor quality of internet-associated false photographs, our investigation showed that deepfake images enhanced using techniques like DFDNet seemed more realistic. When a person is looking directly into the camera, it's difficult to tell that they're in a deepfake. However, flaws in the deepfake-generated picture become more apparent when the individual turns to gaze to their side. We think that datasets with challenging settings like these should be the main emphasis for improving deepfake detection performance. We want to use picture improvement techniques to create deepfakes with fewer flaws and better quality in the future, with the goal of making them more difficult to identify using deepfake detection methods.

## References

- [1]. Wiley, Victor and Lucas, Thomas. (2018). Computer Vision and Image Processing: A Paper Review. International Journal of Artificial Intelligence Research.
- [2]. Bharathi, S.Shankar and N.Radhakrishnan, and Prasad, Pinnamaneni. (2013). Machine Vision Solutions in Automotive Industry.
- [3]. M. H. Wagdy, H. A. Khalil and S. A. Maged, "Swarm Robotics Pattern Formation Algorithms," 2020 8th International Conference on Control, Mechatronics and Automation (ICCMA), 2020, pp. 12-17.
- [4]. Badrinarayanan, V., Kendall, A., and Cipolla, R. (2017). SegNet: A deep convolutional encoder-decoder architecture for image segmentation. IEEE Transactions on Pattern Analysis and Machine Intelligence, 39(12), 2481-2495.
- [5]. Yang, W., Hui, C., Chen, Z., Xue, J. H., and Liao, Q. (2019). FV-GAN: Finger vein representation using generative adversarial networks. IEEE Transactions on Information Forensics and Security, 14(9), 2512-2524.
- [6]. TensorFlow. Accessed on: December 19, 2020. Available at <https://www.tensorflow.org/>
- [7]. AKaliyar, R. K., Goswami, A., and Narang, P. (2020). Deepfake: improving fake news detection using tensor decomposition based deep neural network. Journal of Supercomputing.
- [8]. Deepfake Detection Challenge Results. Accessed on: January 15, 2021. Available at <https://ai.facebook.com/blog/deepfake-detection-challenge-results-an-open-initiative-to-advance-ai/>
- [9]. Dolhansky, Brian and Howes, Russ and Pflaum, Ben and Baram, Nicole and Ferrer, Cristian.(2019). The Deepfake Detection Challenge(DFDC)Preview Dataset.
- [10]. Faceswap: Deepfakes software for all. Accessed on: February 2, 2021. Available at <https://github.com/deepfakes/faceswap>
- [11]. FakeApp 2.2.0. Accessed on: February 2, 2021. Available at <https://www.malavida.com/en/soft/fakeapp/>
- [12]. Balle, Johannes and Laparra, Valero and Simoncelli, Eero.(2016). End-to-end Optimized Image Compression.
- [13]. Cheng, Zhengxue and Sun, Heming and Takeuchi, Masaru and Katto, Jiro. (2019). Energy Compaction-Based Image Compression Using Convolutional AutoEncoder. IEEE Transactions on Multimedia. PP. 1-1.
- [14]. A. Punnappurath and M. S. Brown, "Learning Raw Image Reconstruction-Aware Deep Image Compressors," in IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 42, no. 4, pp. 1013-1019, 1 April 2020.
- [15]. Petrov, Ivan, et al. (2020). DeepFaceLab: A simple, flexible and extensible face swapping framework. arXiv preprint arXiv:2005.05535, 2020.
- [16]. Tewari, Ayush, et al. (2018). High-Fidelity Monocular Face Reconstruction Based on an Unsupervised Model-Based Face Autoencoder. IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [17]. Deep Fakes, Fake News, and What Comes Next. Accessed on: January 20, 2021. Available at <https://jsis.washington.edu/news/deep-fakes-fake-news-and-what-comes-next/>
- [18]. Korshunova, I., Shi, W., Dambre, J., and Theis, L. (2017). Fast faceswap using convolutional

neural networks. In Proceedings of the IEEE International Conference on Computer Vision (pp. 3677-3685).

- [19]. Mirsky, Yisroel and W. Lee. "The Creation and Detection of Deepfakes." ACM Computing Surveys (CSUR). vol. 54 no. 1 , April 2021.
- [20]. D. Afchar, V. Nozick, J. Yamagishi and I. Echizen, "MesoNet: a Compact Facial Video Forgery Detection Network," 2018 IEEE International Workshop on Information Forensics and Security (WIFS), 2018, pp. 1-7.
- [21]. V. Schetinger, M. M. Oliveira, R. da Silva, and T. J. Carvalho. Humans are easily fooled by digital images. arXiv preprint arXiv:1509.05301, 2015.