





Deepfakes are digitally generated media where one person's image is swapped with another's. While this technology has valid uses, like producing convincing visual effects for entertainment, it has also been employed to generate false news and spread misleading information (Patel and Jain, 2024). These manipulated images and videos can realistically change faces, expressions, or even create entirely false identities. This poses serious threats to credibility, personal privacy, security, political stability, and public trust. Because these sophisticated forgeries are becoming easier to accessible and more convincing, it is now critically important to develop good techniques/ strategies to detect them. (Abdullah et al., 2024).

Among various detection strategies, deep learning approaches especially convolutional neural networks (CNNs) have shown remarkable promises to learn fine-grained visual features, detect inconsistencies typical in synthetic images and scale well with data and generalize across manipulation types (Kroiß & Reschke, 2025). One of the most frequently used CNN architectures in image analysis tasks is ResNet-50, a 50-layer residual network known for its balance of depth and computational efficiency. ResNet-50 has been employed across numerous studies for the detection of manipulated media, leveraging its ability to learn variations between real and fake images (Kroiß & Reschke, 2025).

However, despite growing research in this area, there is a lack of consolidated understanding regarding the effectiveness of ResNet-50 in the context of deepfake image detection. Studies differ in terms of datasets used, preprocessing techniques, evaluation metrics, and types of manipulations considered. These differences make it difficult to draw clear conclusions about the model's generalizability and real-world applicability. To address this gap, this study presents a systematic literature review (SLR) focused on deepfake image detection using ResNet-50. The objective is to evaluate and synthesize existing research to identify trends, challenges, and best practices.

By systematically analyzing peer-reviewed studies, our contributions are to provide a comprehensive overview of:

- Deepfake generated methods
- Deepfake image detection methods
- CNN models for deepfake image detection
- The performance of ResNet-50 in various deepfake image detection scenarios,
- The types of datasets, evaluation metrics used and performance compared to other model,
- Limitations, open issues, and future directions.

This review serves as a valuable resource for researchers, practitioners, and policymakers interested in advancing the development of robust, interpretable, and scalable deepfake detection systems. The remainder of the paper is organized as follows: Section II presents the review procedure by defining interest research questions. In Section III, we thoroughly discuss the findings from different studies. Finally, Section VI concludes the paper.

## 2. PROCESS OF SLR

We utilize their methods in our systematic literature review (SLR) and organize the review procedure into three key phases, depicted in Figure 1. This is done to identify, evaluate, and understand various research related to particular research questions.

- **Planning the Review:** This initial phase aims to (a) determine the necessity for the review, (b) establish the standards and methods for conducting it, and (c) assess the suitability of those standards and methods for this specific systematic literature review.
- **Conducting the Review:** Following the guiding principles outlined in (Rana et al., 2022), this phase comprises six essential steps.
  - i. **Research Questions (RQs):** The purpose of defining research questions is to pinpoint the relevant studies that should be included in the current review. We establish a series of research questions (detailed later) specifically focusing on the field of Deepfakes.
  - ii. **Search strategy (SS):** A planned search approach is used to identify a broad range of relevant studies related to our research questions. We aim to develop an objective and comprehensive search strategy to capture as much pertinent literature as possible.
  - iii. **Study Selection Criteria (SSC):** Selecting relevant literature can be challenging, due to potential biases such as the language of publication, familiarity with the authors, institutions, journals, or publication year (Rana et al., 2022). Before defining selection criteria, we carefully consider factors to ensure objectivity in choosing primary studies that offer meaningful evidence related to our research questions.
  - iv. **Quality Assessment Criteria (QAC):** Assessing the quality of each selected study is essential to ensure the findings are reliable and impartial. We develop a set of quality criteria for evaluating the individual studies.
  - v. **Data extraction and monitoring (DEM):** We carefully plan how to collect the necessary data from the chosen studies and document the corresponding evidence.
  - vi. **Data Synthesis (DS):** Data synthesis involves organizing and summarizing the findings from the selected studies. We use a structured set of methods to effectively integrate the information.

**Reporting the Review.** After completing the review of all the studies, we present the results in an appropriate format for distribution and the intended audience.

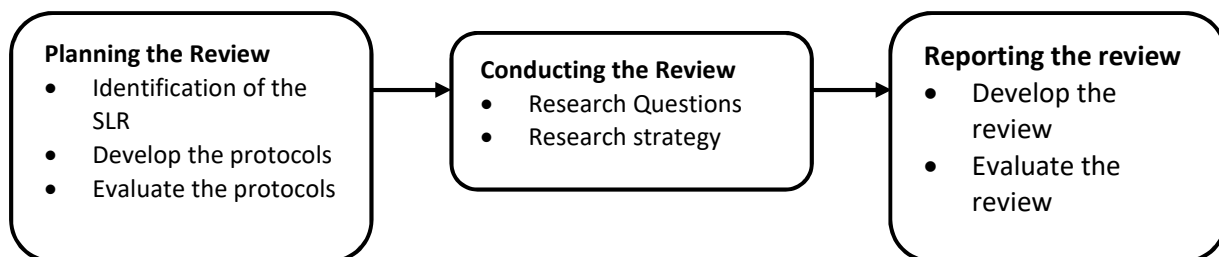


Fig. 1: The process of SLR









### Inclusion and exclusion criteria

We considered different inclusion and exclusion criteria to gather relevant studies for our systematic review. The criteria used to select articles were as follows: **Inclusion criteria** (i) published articles/journals/conference papers between 2020 to 2025 (ii) articles that proposed and developed deepfake image detection CNN (ResNet-50), methods, effectiveness of ResNet-50, modified ResNet-50, dataset and metrics to identify the deepfake image detection; and (iii) articles related to research questions and studies that focus on deepfakes image detection using ResNet-50 only. **Exclusion criteria** articles were as follows: (i) duplicate articles found via different databases; (ii) articles written in other languages; (iii) articles that do not consider deepfakes image detection using ResNet-50 (iv) irrelevant studies (video only or audio only deepfake studies) and (v) survey articles. Using this strategy, we ensured that all related studies were included and that inapplicable studies were excluded from the system article literature review.

### Quality Evaluation

Assessing the strength of the evidence in a systematic review is just as crucial as analyzing the information itself. Problems in research design can skew the results of weaker studies, and this should be acknowledged carefully. Studies with flawed methods or limitations that could affect the accuracy of the results should be identified in, or removed from, literature reviews. Choosing the right criteria to assess the trustworthiness of information and possible problems in each study is also essential. Following established standards like those in the PICO Portal (Pico, 2022), we used such criteria to verify the chosen studies and to examine them for relevant information in this systematic review.

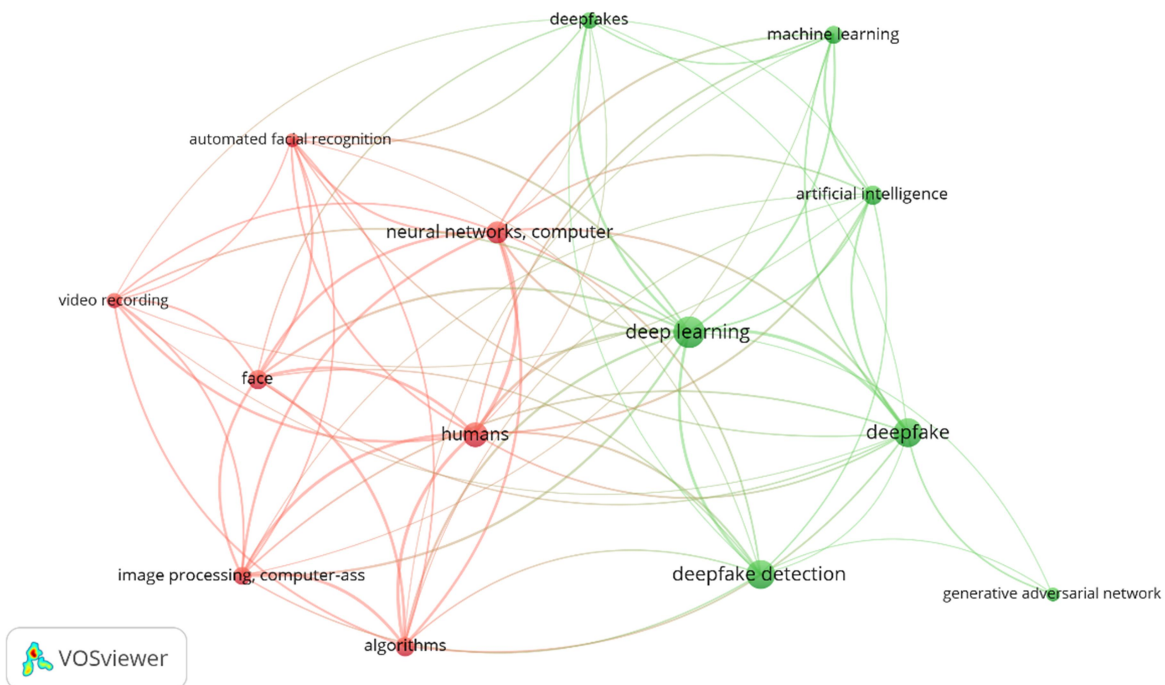


Fig. 5. Keyword co-occurrences visualization analysis of eligible studies using the VOSviewer tool.

Moreover, we employed a confirmation process to judge the articles included, ensuring consistency across different sets of information. After evaluating the quality, we found 45 relevant studies discussing deepfake image detection using ResNet-50. Figure 3 shows the analysis of how frequently different keywords appear together in the relevant studies for our systematic literature review.

### 3. DEEPAKE GENERATION METHODS

Deep learning generation methods have dramatically transformed the deepfake landscape, enabling the creation of highly convincing fabricated media, techniques like Autoencoders, VAE, GAN, Diffusion models (Heo et al., 2021). Fig. 2 illustrates taxonomy of deepfake generation.

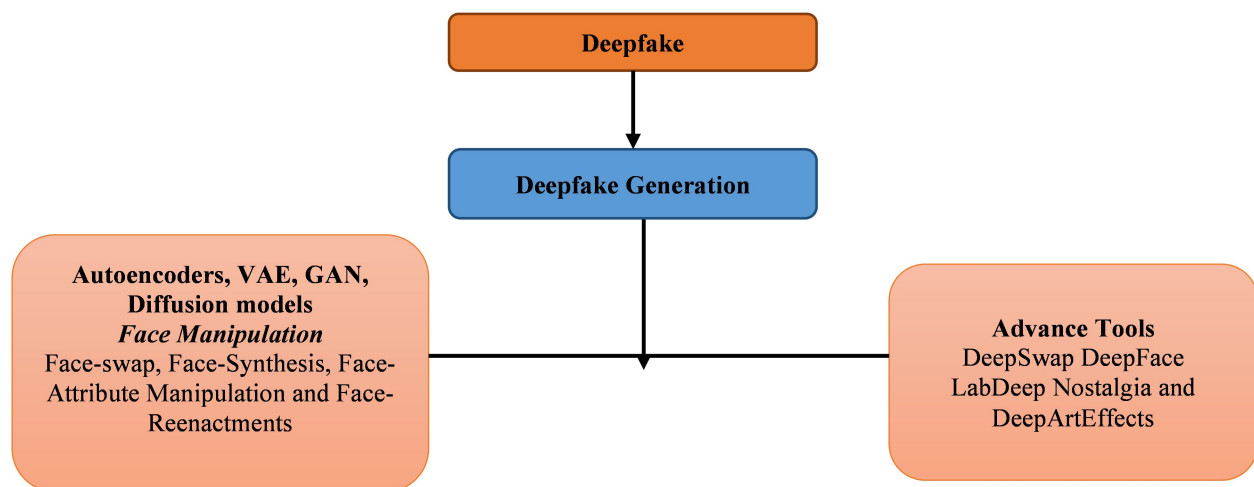


Fig. 6. Taxonomy of Deepfake Generation Methods

**Autoencoders:** It consists of two main components: an encoder and a decoder. The encoder analyzes digital media to identify key characteristics and compress them into a simplified representation (a code that removes irrelevant details). The decoder then uses this representation to reconstruct the original media. By training the encoder and decoder on a varied dataset of authentic and manipulated videos, the autoencoder learns to apply the features of one image to another. This makes it useful for face replacement. However, a significant limitation is that the generated image may be easily detected by humans as artificial. (Kaushik et al., 2025).

**Generative Adversarial Networks:** They tools are made of two parts– the Generator (G) and Discriminator (D). The generator is used to produce novel data samples (fake data), and the discriminator evaluates them against the existing data to draw comparison and conclude whether the content is fake or real. Generative Adversarial Network (GAN) techniques for generating paraphrases are often centered around facial imagery.



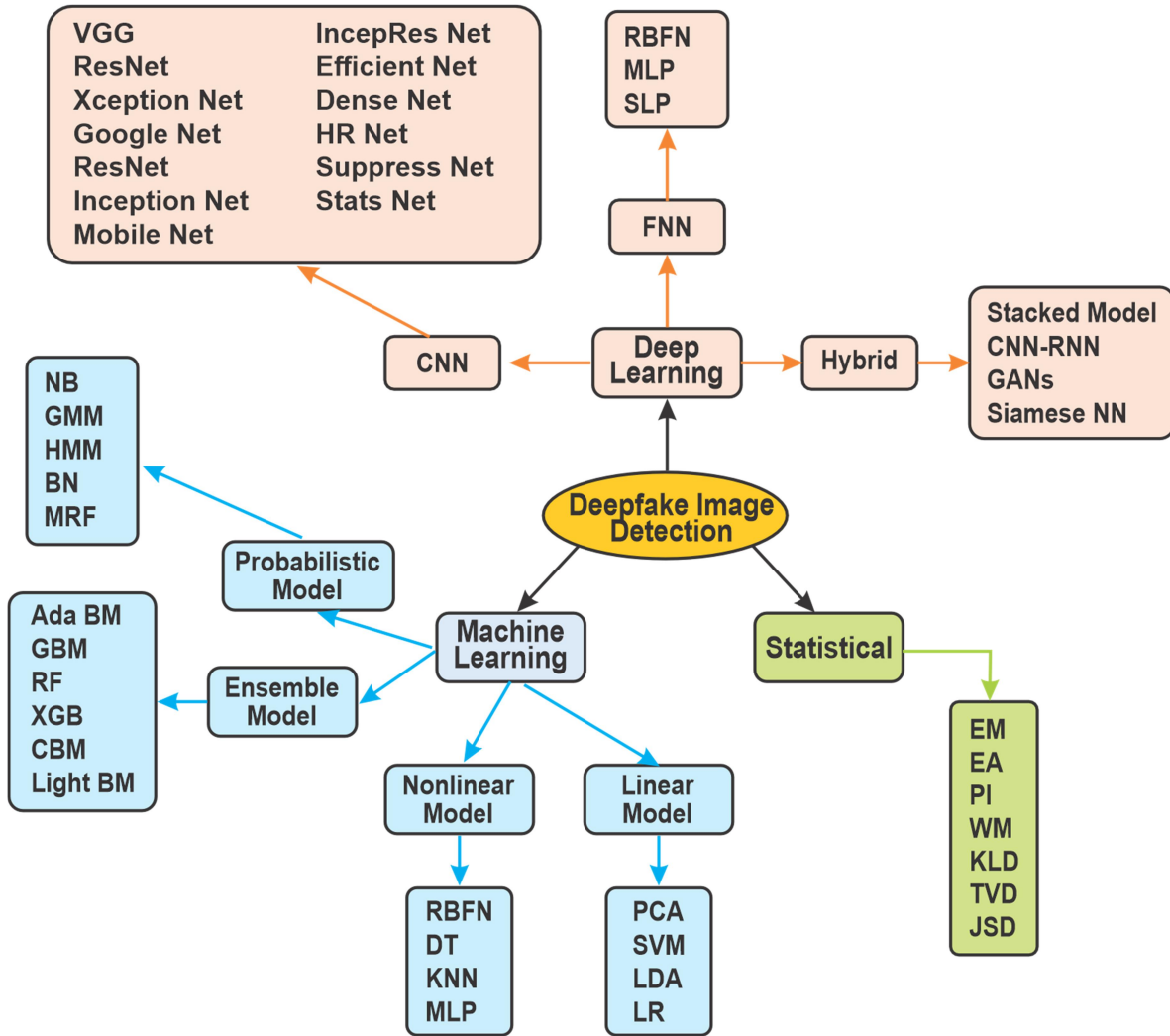


Fig. 7. Taxonomy of major Deepfake image detection methods

### Machine Learning Method

Traditional machine learning methods are beneficial for interpreting the reasoning behind any outcome that can be expressed in understandable terms. Because machine learning approaches address many current challenges, they are well-suited for tasks like deepfake detection, offering improved data and process understanding. Additionally, adjusting settings and modifying model designs is much simpler. Generative Adversarial Networks (GANs), on the other hand, are self-training generative models that handle unsupervised problems as supervised, producing highly realistic fake faces in images or videos. (Agarwal et al., 2021) explores deepfake detection using Support Vector Machines (SVM).





By using techniques to make the model resistant to adversarial attacks and extracting comprehensive features, the approach effectively identifies subtle traces left behind during image tampering. The findings demonstrate that the model can be applied to various types of image manipulation and different datasets, making it a practical and dependable solution for real-world use. The study highlights the value of combining different detection methods to handle the intricate nature of fake media analysis and emphasizes the need to integrate features from multiple domains to tackle the increasing challenges posed by deepfakes. (Gura et al., 2024) proposed a novel Convolutional Neural Network (CNN) approach for identifying deepfake videos. The method begins by extracting key facial features from video frames using facial landmark detection. This structured facial data is then fed into a specialized CNN. To improve the model's performance, a data augmentation technique is employed within the CNN itself, creating synthetic "fake data" or "fake images" to train on.

The system was developed using Python and associated libraries. The researchers used footage from the Deep Fake Detection Challenge dataset, specifically 242 videos (199 fake, 53 real) and an additional 76 (40 fake, 36 real), each segment being 10 seconds long. The proposed CNN achieved a high testing accuracy of 91.47%, with a loss of 0.342 and an AUC score of 0.92. These results surpassed the performance of two comparative methods: a standard CNN and an MLP-CNN. Notably, this new method also demonstrated superior accuracy compared to existing state-of-the-art models, including XceptionNet, Meso-4, EfficientNet-BO, Mesoinception-4, VGG-16, and DST-Net. The core innovation of this work lies in creating a new CNN-based learning model specifically designed for precise detection of manipulated facial images in deepfakes. (Kosarkar et al 2023) identify deepfake images within a video dataset, we employed a specifically designed convolutional neural network (CNN) algorithm. We then compared its performance against two alternative approaches to ascertain the optimal method. Our model was trained and tested using a dataset from Kaggle.

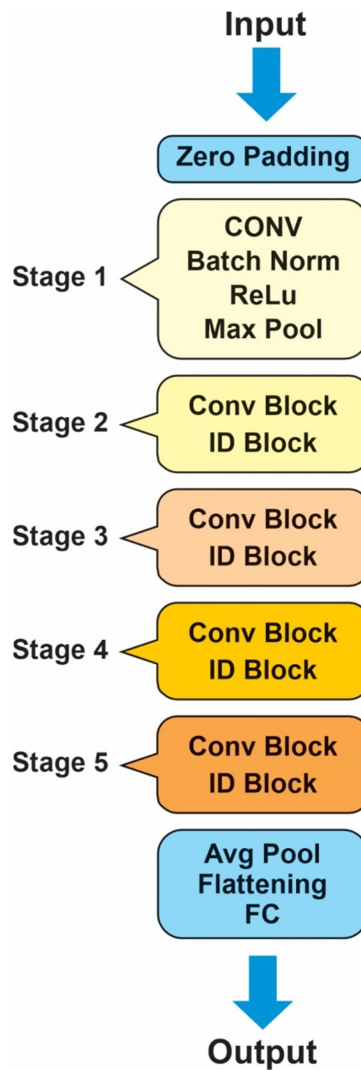
This research leverages CNNs to differentiate between genuine and manipulated images by training three separate CNN models. Furthermore, we created and implemented a refined CNN model, incorporating additional components such as a fully connected layer, maximum pooling, and a dropout layer. The methodology involves extracting frames, identifying facial features, preparing the data, and finally classifying the images as either authentic or fabricated, aligning with the research aims. Performance was evaluated using metrics like accuracy, loss, and the area under the ROC curve. Our enhanced CNN model demonstrated superior results, achieving 91.4% accuracy, a lower loss value of 0.342, and an AUC of 0.92. In addition, the standard CNN model reached 85.2% accuracy on the test set, while the MLP-CNN model achieved 95.5% testing accuracy.

### **Statistical Measurements Methods**

Statistical methods have become important resources for identifying deepfakes. These methods use complex statistical analysis to find trends and differences in data, helping to spot altered content. By studying various statistical characteristics, such as image quality measures, noise variations, facial points, and timing problems, these approaches aim to distinguish between real and fake media. (Hou et al., 2023) proposed a statistically-driven attack, termed Stat Attack, has been introduced to challenge DeepFake detectors. This approach operates in two stages. Initially, a collection of naturally occurring distortions (specifically, exposure, blur, and noise) are deliberately introduced to the synthesized images to undermine the detector.







**Fig. 8 Taxonomy of ResNet-50 Architecture**

(Singh & Ramachandra, 2022) presents a Deep Fake Detection system (DFD) powered by deep learning. The system utilizes enhanced convolutional layers before feeding into a Resnet-50 architecture. DFD is trained from start to finish using low-quality images from the FaceForensics++ dataset. Around 1.68 million images are used for training, 315,000 for validation, and 340,000 for testing. The DFD is evaluated in three settings: a mix of various image manipulations (achieving 96.07% accuracy, significantly better than the existing best methods at 85.14%), individual manipulation techniques (achieving perfect 100% accuracy on neural textures), and manipulations involving different images (achieving 94.28% accuracy on the previously unseen faceswap category, outperforming the state-of-the-art).





















			precision, ROC, F1-score AUC	(face analysis, synthesis, face attribute and classification)
(Berrahal et al., 2023)	CelebA-HQ	Image	AUC	face attributes/identity, deepfake and spoof detection using ResNet-50 as binary classifier
(Chouhan et al., 2024)	CELEB-DF	Image/Video	Accuracy	Deepfake video detection/facial forgery classification used as binary classifier (real or fake)
(Safwat et al., 2024)	FF++	Image/Video	Accuracy	Deepfake image detection, train ResNet-50 for fake or real face classification
(Samal et al., 2024)	Open Forensics	Image	Accuracy	Deepfake/image manipulation detection, evaluating generalization to unseen manipulation
(Patel et al. 2023)	AttGAN, Image STAR GAN, StyleGAN and StyleGAN2		Accuracy	face attribute manipulation, deepfake images editing classify manipulated or original faces

### RQ 3: What evaluation metric are used to computing the performance of deepfake image detection model using ResNet-50?

This section outlines several evaluation metrics used to assess how well models perform in identifying manipulated images and videos. A confusion matrix summarizes the results of the classification, showing both correct and incorrect predictions. The effectiveness of the methods is evaluated and validated using data from this matrix. Specifically, True Positives (TP) represent the number of Deepfakes correctly identified as Deepfakes, while True Negatives (TN) indicate the









## REFERENCES

1. Abdal Mashkour, K. A. (2025). Classifying Real and AI-Generated Images Using Fine-Tuned ResNet50. *Journal of Advanced Artificial Intelligence, Engineering and Technology*, 01(02). <https://doi.org/10.56147/aaiet.1.2.7>
2. Abdullah, S. M., Cheruvu, A., Kanchi, S., Chung, T., Gao, P., Jadliwala, M., & Viswanath, B. (2024). An Analysis of Recent Advances in Deepfake Image Detection in an Evolving Threat Landscape (No. arXiv:2404.16212). arXiv. <https://doi.org/10.48550/arXiv.2404.16212>
3. Agarwal, S., Girdhar, N., & Raghav, H. (2021, January 28). A Novel Neural Model based Framework for Detection of GAN Generated Fake Images. 2021 11th International Conference on Cloud Computing, Data Science & Engineering (Confluence). 2021 11th International Conference on Cloud Computing, Data Science & Engineering (Confluence), Noida, India. <https://doi.org/10.1109/confluence51648.2021.9377150>
4. Aghasanli, A., Kangin, D., & Angelov, P. (2023). Interpretable-through-prototypes deepfake detection for diffusion models. 2023 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW), 467–474. <https://doi.org/10.1109/iccvw60793.2023.00053>
5. Akram, O., Mohamed, A., Magdy, H., Abdellatif, M.M., Abdelghafar, S. (2025). Comparative Analysis of Custom CNN Architecture and MobileNet for Deepfake Image Detection. *Proceedings of the 11th International Conference on Advanced Intelligent Systems and vol 238*. [https://doi.org/10.1007/978-3-031-81308-5\\_6](https://doi.org/10.1007/978-3-031-81308-5_6)
6. Alzahrani, A. (2024). Digital Image Forensics: An Improved DenseNet Architecture for Forged Image Detection. *Engineering, Technology & Applied Science Research*, 14(2), 13671–13680. <https://doi.org/10.48084/etasr.7029>.
7. Alzamily, J. Y. I., Ariffin, S. B., & Naser, S. S. A. (2022). CLASSIFICATION OF ENCRYPTED IMAGES USING DEEP LEARNING –RESNET50. . . Vol., 21.
8. Amrithat, D.V & Philimon, S., (2025): Deep GuardNet: A Novel CNN Architecture for Deep Fake Image Detection. *International Conference on Machine Learning and Data Engineering Procedia Computer Science*. 258, 811–818.
9. Bartos, G. E., & Akyol, S. (2024). Deep Learning for Image Authentication: A Comparative Study on Real and AI-Generated Image Classification.
10. Borade, S., Jain, N., Patel, B., Kumar, V., Godhrawala, M., Kolaskar, S., Nagare, Y., Shah, P., & Shah, J. (2024). ResNet50 DeepFake Detector: Unmasking Reality. *Indian Journal of Science and Technology*, 17(13), 1263–1271. <https://doi.org/10.17485/ijst/v17i13.285>
11. Berrahal, M., Boukabous, M., Yandouzi, M., Grari, M., & Idrissi, I. (2023). Investigating the effectiveness of deep learning approaches for deep fake detection. *Bulletin of Electrical Engineering and Informatics*, 12(6), 3853–3860. <https://doi.org/10.11591/eei.v12i6.6221>
12. Bird, J. J., & Lotfi, A. (2024). CIFAKE: Image Classification and Explainable Identification of AI-Generated Synthetic Images. *IEEE Access*, 12, 15642–15650. <https://doi.org/10.1109/access.2024.3356122>
13. Brian, D., Joanna, B., Ben, P., Jikuo, Lu., Russ, H., Menglin, W., & Cristian, C. F. (2024): Dataset: Deepfake Detection Challenge (DFDC). <https://doi.org/10.57702/rdr5mk24>
14. Cao, X., & Gong, N. Z. (2021). Understanding the Security of Deepfake Detection (No. arXiv:2107.02045). arXiv. <https://doi.org/10.48550/arXiv.2107.02045>







