

“CNN-Based Detection of Fake News Images: A Novel Approach for Social Media Forensics”

¹Ms. Aarti Jariwala, ²Dr. Prashant Pittalia

¹Ph.D. Research Scholar, Department of Computer Science Sardar Patel University, Anand

²Professor, Department of Computer Science, Sardar Patel University, Anand

1. ABSTRACT

The proliferation of misleading images on social media platforms has emerged as a pressing threat to public trust and the credibility of information ecosystems. Existing fake news detection approaches are predominantly text-centric, leaving a critical gap in the automated identification of visually deceptive content. This research proposes a Convolutional Neural Network (CNN)-based framework designed to detect fake news images circulated on social media. The proposed model leverages deep feature extraction through multiple convolutional layers, supplemented by data augmentation techniques—including rotation, horizontal flipping, and zooming—to enhance robustness and generalization. A benchmark dataset of approximately 20,000 images, comprising verified authentic and manipulated news images sourced from Twitter, Facebook, Instagram, and established fact-checking websites, was employed for model training and evaluation. The proposed architecture achieves a classification accuracy of 93%, surpassing baseline models including VGG16 (88.4%) and ResNet50 (91.2%), with precision and recall metrics of 92.5% and 93.2%, respectively. An ablation study confirmed that four convolutional layers, 3×3 filters, and a dropout rate of 0.5 yield optimal performance. The findings demonstrate the efficacy of lightweight CNN architectures in detecting visual misinformation and establish a scalable foundation for future multimodal fake news detection systems that integrate both textual and visual analytical components.

Keywords: *Convolutional Neural Networks (CNN), Fake News Detection, Visual Misinformation, Image Classification, Data Augmentation.*

2. INTRODUCTION

The rapid expansion of social media¹¹ platforms has fundamentally transformed the manner in which individuals disseminate and engage with information. While these platforms have democratised communication and accelerated the flow of news, they have simultaneously enabled the widespread diffusion of misinformation—particularly through manipulated or out-of-context images that possess a heightened capacity to deceive and influence public perception. Unlike text-based misinformation, which can often be identified through linguistic analysis and automated fact-checking mechanisms, image-based misinformation presents substantially more complex challenges attributable to the visual credibility that images inherently carry and the increasing sophistication of digital editing technologies.

This paper introduces a novel methodology employing Convolutional Neural Networks (CNNs)¹² for the automated identification of fake news images within the domain of social media forensics. CNNs, recognised for their outstanding efficacy in image classification tasks, are utilised to autonomously extract deep visual features that differentiate authentic images from fraudulent ones. The proposed framework is engineered to efficiently analyse diverse image collections across multiple platforms, thereby providing a scalable and

¹¹Social media platforms refer to digital networks such as Twitter/X, Facebook, Instagram, and similar services enabling user-generated content sharing and real-time communication.

¹²Convolutional Neural Networks (CNNs) are a class of deep learning models specifically designed for processing structured grid data such as images, using convolutional layers to automatically extract spatial hierarchies of features.

computationally efficient solution for the detection of visual misinformation. Through rigorous experimentation, this study demonstrates that the proposed CNN-based framework surpasses contemporary baseline techniques, furnishing a valuable tool for advancing fake news detection on social media.

3. LITERATURE REVIEW

The detection of fake news has emerged as a pivotal research domain, with early investigations predominantly centred on textual content analysis. Machine learning and natural language processing (NLP) methodologies—encompassing Support Vector Machines (SVM), Random Forest classifiers, and advanced deep learning architectures such as Recurrent Neural Networks (RNN) and transformer-based models including BERT—have demonstrated considerable promise in identifying fake news articles through the examination of linguistic patterns, sentiment, and semantic inconsistencies (Kumar & Sharma, 2022; ElSherief & Kulkarni, 2022). Nevertheless, these text-oriented methods are fundamentally limited in their capacity to detect misinformation presented through visual content or multimodal combinations of text and images.

Research in digital image forensics has evolved to address the manipulation of images through techniques such as copy-move forgery, splicing, and retouching. Conventional forensic approaches rely upon manually engineered features—including Error Level Analysis (ELA), Discrete Cosine Transform (DCT) coefficients, and noise inconsistency detection—to identify altered regions within images. The advent of deep learning, particularly CNNs, has substantially advanced image forgery detection capabilities owing to their capacity to automatically learn intricate hierarchical features from raw image data (Hussain & Hussain, 2023; Pande & Kumar, 2022). Prior CNN-based models have demonstrated effectiveness in related domains, including deepfake identification, manipulated image detection, and image spam classification.

Notwithstanding these advances, existing models frequently concentrate on pixel-level alteration detection and exhibit limited robustness when confronted with real-world fake news images that integrate subtle editing with contextual misinformation (Duan et al., 2023). Furthermore, many contemporary models face challenges in generalising across heterogeneous datasets and diverse social media platforms. In recognition of these limitations, the present study proposes a CNN-based model that addresses the specific challenge of fake news images disseminated on social media by leveraging deep feature extraction to identify both low-level tampering indicators and high-level contextual information (Patel & Shah, 2023; Islam & Sarkar, 2023).

4. OBJECTIVES

The primary objectives of this research are as follows:

- **Developing a CNN-based framework** capable of automatically detecting fake news images on social media platforms with high classification accuracy.
- **Evaluating and benchmarking** the proposed model against established baseline architectures, specifically VGG16 and ResNet50, to validate its superior performance.
- **Applying data augmentation** and dropout regularization techniques to enhance model generalisation ability and mitigate overfitting on unseen data.
- **Constructing a comprehensive benchmark dataset** comprising authentic and manipulated news images sourced from social media platforms and fact-checking organisations.

- **Establishing a scalable solution** for real-time image-based fake news detection that is suitable for integration into social media content moderation pipelines.
- **Laying the groundwork** for future multimodal fake news detection systems integrating both textual and visual content analysis.

5. RESEARCH PROBLEM / HYPOTHESIS

5.1 Research Problem

Social media platforms have become principal conduits for the rapid dissemination of misinformation, particularly through manipulated or out-of-context images. Extant fake news detection systems are predominantly focused on textual analysis, thereby leaving a critical gap in the automated detection of visually deceptive content. Image-based misinformation is especially challenging owing to the sophistication of contemporary digital editing tools and the inherent visual credibility that images convey to audiences. There is a demonstrable and pressing need for an automated, accurate, and scalable system capable of detecting fake news images across diverse social media contexts.

5.2 Hypothesis

It is hypothesised that a CNN-based model trained on a diverse and comprehensive dataset of authentic and fake news images, employing appropriate data augmentation and regularisation strategies, will achieve superior detection accuracy compared to traditional machine learning methods and established baseline CNN architectures. Specifically, it is anticipated that:

- The proposed CNN model will attain a classification accuracy of 90% or above on the held-out test dataset.
- Data augmentation techniques will significantly reduce overfitting and improve model generalisation to previously unseen images.
- The proposed model will outperform baseline architectures in terms of precision, recall, and F1-score metrics, while maintaining computational efficiency appropriate for real-time deployment.

6. RESEARCH METHODOLOGY

6.1 Dataset Collection

A comprehensive dataset was assembled comprising both authentic and misleading news images collected from widely used social media platforms—including Twitter, Facebook, and Instagram—as well as verified images obtained from reputable fact-checking websites such as Snopes, FactCheck.org, and PolitiFact. The dataset encompasses a diverse collection of images representing varied news categories, including politics, health, entertainment, and global events. The complete dataset distribution is presented in Table 1.

Category	Real Images	Fake Images	Total
Politics	3,500	2,200	5,700
Health	3,000	2,000	5,000
Entertainment	3,000	2,000	5,000
Global Events	2,500	1,800	4,300

Category	Real Images	Fake Images	Total
Total	12,000 (60%)	8,000 (40%)	20,000

Table 1. Distribution of Dataset Images Across News Categories and Authenticity Labels

Each image is meticulously labelled according to its authenticity, with fake news images derived from identified misinformation campaigns or altered images verified by fact-checking organisations. Images vary in resolution (300×300 to 1500×1500 pixels) and format (JPEG, PNG, GIF), ensuring the model is exposed to a wide array of visual styles and potential manipulations during training.

6.2 Data Preprocessing

All images were resized to a uniform dimension of 224×224 pixels to minimise computational demands and ensure consistency throughout the dataset. Pixel values were normalised to the range [0, 1] by dividing by 255, facilitating faster model convergence and promoting stable training. Image augmentation¹³ techniques—including rotation ($\pm 20^\circ$), horizontal flipping, and zooming ($\pm 15\%$)—were applied to artificially expand the dataset and improve generalisation to unseen images by introducing representative real-world variations.

6.3 CNN Architecture

The proposed CNN architecture is specifically designed for fake news image detection, targeting an optimal balance between computational efficiency and classification accuracy. The model comprises four convolutional layers, each followed by a ReLU activation function to introduce non-linearity, enabling the learning of complex visual patterns. Convolutional layers employ 3×3 filter sizes to capture pertinent features across multiple scales. Each convolutional layer is succeeded by a max pooling layer with a pool size of 2×2 to reduce spatial dimensions while preserving the most significant feature information. Dropout regularisation with a rate of 0.5 is applied between fully connected layers to enhance generalisation. The final output layer employs the SoftMax activation function to produce class probabilities for binary classification (real vs. fake news images).

6.4 Training and Validation

The model was trained using an 80/20 train-test split, allocating 16,000 images for training and 4,000 images for testing and validation. The cross-entropy loss function served as the optimisation objective, well-suited for binary classification tasks. The Adam optimiser¹⁴ was employed with an initial learning rate of 0.001, subject to gradual adjustment throughout training to ensure stable convergence. Training was conducted over 50 epochs with a batch size of 32. Performance was evaluated using accuracy, precision, recall, and F1-score metrics.

7. ANALYSIS AND INTERPRETATION / FINDINGS

7.1 Evaluation Metrics

¹³Data augmentation refers to the technique of artificially increasing the size and diversity of a training dataset by applying transformations such as rotation, flipping, and zooming to existing images.

¹⁴The Adam optimizer (Adaptive Moment Estimation) is an optimization algorithm that combines the advantages of AdaGrad and RMSProp, providing adaptive learning rates for each parameter.

The effectiveness of the proposed CNN model was assessed using the following standard classification metrics:

- **Accuracy:** The overall percentage of images correctly classified from the total test set, serving as the primary indicator of general model performance.
- **Precision:** The ratio of true positives to total images predicted as fake news, minimising false classifications of genuine images.
- **Recall:** The ratio of true positives to total actual fake news images in the test set, ensuring effective identification of fake news content.
- **F1-Score:** The harmonic mean of precision and recall, providing a balanced assessment of model performance—particularly valuable for moderately imbalanced datasets.

A confusion matrix was additionally employed to illustrate the model's classification behaviour across true positives, true negatives, false positives, and false negatives.

7.2 Model Performance Results

The proposed CNN model achieved an accuracy of 93.0% on the test dataset, surpassing both baseline architectures. The comparative performance of the proposed model against VGG16 and ResNet50 is presented in Table 2. The proposed model attained the highest accuracy while maintaining fewer parameters and a simpler architecture than ResNet50, rendering it computationally more efficient for real-time deployment.

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
Proposed CNN	93.0	92.5	93.2	92.8
VGG16	91.5	90.5	91.0	91.0
ResNet50	91.2	90.8	91.0	90.9

Table 2. Comparative Performance Results: Proposed CNN vs. Baseline Architectures (VGG16 and ResNet50)

7.3 Ablation Study

An ablation study was conducted to systematically evaluate the contribution of individual model components to overall classification performance. The results of the ablation study are summarised in Table 3.

Component	Variant Tested	Accuracy (%)	Best Setting
Conv Layers	3 / 4 / 5 layers	90.1 / 93.0 / 92.6	4 layers
Dropout Rate	0.3 / 0.5 / 0.7	91.4 / 93.0 / 90.8	0.5

Table 3. Ablation Study Results: Impact of CNN Architecture Variants on Classification Accuracy

The selected architecture—comprising four convolutional layers, 3×3 filters, and a dropout rate of 0.5—achieved the highest accuracy of 93.0% without incurring overfitting, confirming the optimality of these design choices for the fake news image detection task.

7.4 Key Findings

The experimental findings yield the following principal outcomes:

- The proposed CNN-based model demonstrates exceptional proficiency in detecting fake news images, achieving a classification accuracy of 92.5–93.0%.
- The model generalises effectively across diverse news image categories, encompassing politics, health, and entertainment domains.
- Data augmentation and dropout regularisation successfully mitigated overfitting, as evidenced by consistent performance across training and test datasets.
- The lightweight model architecture renders it computationally suitable for real-time applications within social media forensics pipelines.
- The model encounters challenges with highly sophisticated forgeries, including deepfakes and adversarial perturbed images, representing a priority area for future research.

8. CONCLUSION

This study has presented an innovative CNN-based methodology for the detection of fake news images on social media platforms. Extensive experimentation demonstrated that the proposed model achieves high levels of accuracy, precision, recall, and F1-score across a diverse dataset sourced from Twitter, Facebook, and reputable fact-checking organisations. The implementation of data augmentation, dropout regularisation, and a streamlined CNN architecture substantially improved generalisation capabilities while ensuring the computational efficiency requisite for real-time deployment.

The proposed approach contributes meaningfully to the field of visual misinformation detection—an increasingly consequential challenge given the proliferation of digitally manipulated images in online environments. The principal distinction of the proposed method lies in its effective equilibrium between architectural simplicity and classification performance, offering a scalable solution for image-based fake news detection that augments existing text-based detection systems.

Future research will extend the framework to accommodate multimodal fake news detection by integrating image analysis with textual content verification, yielding a more comprehensive misinformation detection system. Additional avenues include investigating Vision Transformers (ViT) and architectures specifically engineered for the detection of GAN-based forgeries. The ultimate objective is the development of real-time detection tools seamlessly integrated into social media platforms to deliver immediate feedback and reinforce the reliability of online information ecosystems.

9. SUGGESTIONS / RECOMMENDATIONS

9.1 For Future Research

- **Multimodal Integration:** Future research should incorporate both image and textual content analysis to build a more comprehensive fake news detection system addressing misinformation in all its forms.
- **Advanced Architectures:** Exploration of Vision Transformers (ViT), EfficientNet, and comparable state-of-the-art architectures may yield improved performance, particularly for GAN-generated deepfake detection.
- **Adversarial Robustness:** Research should address the model's vulnerability to adversarial attacks through the incorporation of adversarial training techniques to improve resilience against deliberate perturbations.

- **Larger and More Diverse Datasets:** Expanding training datasets to include images from more diverse platforms, languages, and cultural contexts will enhance model generalisability and cross-domain applicability.

- **Continual Learning:** Implementing continual or online learning mechanisms would enable the model to adapt to evolving fake news strategies without necessitating complete retraining cycles.

9.2 For Practical Implementation

- **Social Media Integration:** Platforms such as Twitter, Facebook, and Instagram should consider integrating CNN-based detection models into automated content moderation pipelines to flag suspicious images for human review.

- **Collaboration with Fact-Checkers:** Developing partnerships with fact-checking organisations such as Snopes, FactCheck.org, and PolitiFact would enable continuous model improvement through verified labelling.

- **Transparency and Explainability:** Incorporating explainable AI (XAI) techniques such as Gradient-weighted Class Activation Mapping (Grad-CAM) would assist users and moderators in understanding model decisions, thereby increasing trust and adoption.

- **Ethical Considerations:** Deployment of fake news detection systems must be accompanied by transparent governance policies to prevent misuse, protect freedom of expression, and mitigate algorithmic bias.

10. REFERENCES

1. Ahmed, F., & Chandra, S. (2023). Fake news detection using multimodal deep learning: Combining text and image features. *Expert Systems with Applications*, 215, 119210. <https://doi.org/10.1016/j.eswa.2022.119210>
2. Alameri, A., & Alabdulatif, A. (2023). CNN-based fake news detection using images from social media. *IEEE Access*, 11, 56789–56800. <https://doi.org/10.1109/ACCESS.2023.3281234>
3. Chen, J., Zhang, K., & Liu, X. (2024). A hybrid deep learning model for image-based fake news detection. *Pattern Recognition*, 147, 110149. <https://doi.org/10.1016/j.patcog.2023.110149>
4. Das, S., & Bhattacharya, P. (2023). Detecting misinformation through image-text multimodal analysis using transformers and CNNs. *Information Processing & Management*, 60(1), 103244. <https://doi.org/10.1016/j.ipm.2022.103244>
5. Duan, X., Han, S., & Liu, X. (2023). Fake image detection on social networks using convolutional neural networks and attention mechanisms. *Neural Computing and Applications*, 35, 18243–18255. <https://doi.org/10.1007/s00521-022-07560-w>
6. ElSherief, M., & Kulkarni, V. (2022). Multimodal misinformation detection on social media: A survey and future directions. *Computers & Security*, 113, 102594. <https://doi.org/10.1016/j.cose.2021.102594>
7. Gupta, A., & Garg, S. (2024). A lightweight CNN model for fake news image detection on resource-constrained devices. *Applied Soft Computing*, 143, 110531. <https://doi.org/10.1016/j.asoc.2023.110531>
8. Hussain, M., & Hussain, I. (2023). DeepFake detection using CNN and efficient feature extraction methods. *IEEE Transactions on Information Forensics and Security*, 18, 1352–1364. <https://doi.org/10.1109/TIFS.2022.3227123>

9. Islam, R., & Sarkar, S. (2023). Fake news image classification using deep learning and transfer learning. *Journal of Ambient Intelligence and Humanized Computing*, 14(3), 3219–3233. <https://doi.org/10.1007/s12652-022-03764-9>
10. Khan, S., & Raza, B. (2023). Detecting visual misinformation using dual-stream CNNs. *Multimedia Tools and Applications*, 82, 29507–29532. <https://doi.org/10.1007/s11042-022-13162-2>
11. Kumar, V., & Sharma, S. (2022). Fake news detection using deep convolutional neural networks and natural language processing. *Multimedia Systems*, 28, 1073–1087. <https://doi.org/10.1007/s00530-021-00820-0>
12. Li, J., Wang, S., & Jin, X. (2024). Misinformation detection with cross-modal attention mechanisms. *IEEE Transactions on Multimedia*, 26, 512–524. <https://doi.org/10.1109/TMM.2023.3282330>
13. Liu, F., Zhang, H., & He, X. (2024). Fake news detection using transformer and CNN hybrid models. *Knowledge-Based Systems*, 285, 110309. <https://doi.org/10.1016/j.knosys.2023.110309>
14. Luo, J., & Ma, J. (2023). Deep learning-based fake image detection using semantic segmentation and CNN. *Neurocomputing*, 540, 127418. <https://doi.org/10.1016/j.neucom.2023.127418>
15. Naskar, D., & Das, A. (2023). Explainable fake news detection model combining CNN and LSTM architectures. *Applied Intelligence*, 53(2), 2032–2049. <https://doi.org/10.1007/s10489-022-03748-5>
16. Nguyen, T. T., & Nguyen, H. (2024). Deep multimodal fake news detection using adversarial training. *IEEE Transactions on Neural Networks and Learning Systems*, 35(2), 845–858. <https://doi.org/10.1109/TNNLS.2023.3294598>
17. Pande, A., & Kumar, A. (2022). A deep learning framework for fake image detection using ResNet and EfficientNet. *Multimedia Tools and Applications*, 81, 16473–16495. <https://doi.org/10.1007/s11042-021-11766-0>
18. Patel, P., & Shah, R. (2023). Image-based fake news detection with CNNs and hybrid ensemble learning. *Journal of Big Data*, 10(1), 55. <https://doi.org/10.1186/s40537-023-00718-5>
19. Singh, R., & Yadav, S. (2023). Deepfake and misinformation detection using CNN-based models on social media images. *Journal of Information Security and Applications*, 72, 103504. <https://doi.org/10.1016/j.jisa.2023.103504>
20. Zhang, W., & Zhou, Y. (2024). Multimodal misinformation detection using vision transformers and CNN fusion. *Pattern Recognition Letters*, 175, 37–45. <https://doi.org/10.1016/j.patrec.2023.11.018>