

Document Forgery Detection

Jenifer J¹, Sathiyapriya G², Agalya R³, Supriya V⁴, Shafreen Banu S⁵

^{1, 2, 3, 4, 5}Department of Computer Science and Engineering Jai Shriram Engineering College, Tamil Nadu, India.

Email ID: jenifercse@jayshriram.edu.in¹, gsathyapriya05@gmail.com², agalyaramasamy118@gmail.com³, supriyavelmurugan@gmail.com⁴, shafreenbanu123@gmail.com⁵

Abstract

The rapid growth of digital documentation has increased the risk of document forgery in several industries, including public services, banking, and education. A significant part of traditional document verification methods is manual inspection, which is often inefficient and prone to errors. To detect forged documents, an automated technique that utilizes machine learning and image processing is presented in this paper. The system examines structural and content-based features to identify unauthorized changes, such as altered text areas and visual elements. The proposed model uses preprocessing and feature extraction techniques to differentiate between genuine and fraudulent document types. The results demonstrate that the system improves verification accuracy while decreasing the need for human intervention. This strategy offers a practical means of enhancing the security of documents.

Keywords: Document Forgery Detection, Machine Learning, Image Processing, Automated Verification, Feature Extraction, Fraud Detection, Digital Documents, Document Security, Forged Document Identification, Pattern Recognition

1. Introduction

The rapid advancement of digital technology has led to the widespread use of documents in sectors such as government services, banking, healthcare, and education, among others Figure 1. Important documents that are regularly shared and stored digitally include financial records, identity documents and certificates[1]. This widespread use has also raised the risk of document forgery, which is the alteration or manipulation of documents to obtain unapproved advantages Figure 2. Document forgery includes text alteration, image manipulation, copy-paste forgery, and signature tampering[2]. Traditional document verification methods rely on manual inspection, which is time-consuming, labor-intensive, and prone to human error[3]. As the number of documents increases, manual verification becomes unreliable and ineffective in ensuring data quality[4]. To address these issues, automated document forgery detection systems have garnered significant interest[5]. Image processing and machine learning techniques offer useful tools for analyzing document features and identifying irregularities that indicate forgeries[6]. These techniques can be used to search for difficult-to-find textual, visual, and

structural patterns[7]. This study presents an automated approach that uses machine learning and image processing techniques to detect document forgeries Figure 3. The primary goals of the proposed system are to preprocess documents, extract relevant features, and classify documents as authentic or fraudulent[8]. The goals of this study are to improve accuracy, reduce verification time, and minimize human intervention in document authentication. Digital transformation has resulted in a significant increase in the volume of documents processed electronically in recent years. Identity cards, financial records, contracts, and educational certificates are frequently exchanged in digital formats. This modification increases operational effectiveness, but also exposes systems to risks such as unauthorized modifications and document forgery[9]. The volume of documents processed electronically has dramatically increased in recent years owing to digital transformation. Digital formats are widely used for the exchange of identity cards, financial records, legal contracts and educational certificates Figure 4. Although this change improves operational efficiency, it also leaves systems vulnerable to threats

such as document forgery and unauthorized changes Figure 5. Forgery techniques have developed over time, from straightforward manual changes to complex digital manipulations employing cutting-edge editing tools[10]. Because of these techniques, forged documents appear visually similar to authentic documents, making manual detection challenging. Forged documents are frequently used to obtain illicit advantages, such as jobs, loans, or admission to schools Figure 6. Document forgery results in economic losses, legal disputes, and harm to an organization's reputation. As a result, confirming the authenticity of documents has become essential in recent years. Automated forgery detection systems offer a promising solution by decreasing reliance on manual verification and increasing consistency and dependability. The goal of this study is to create an automated system that uses computational methods to examine document characteristics and spot possible forgeries. The proposed method seeks to assist organizations in making verification decisions more quickly and accurately Figure 7.

2. Related Work

In recent years, document forgery detection has been extensively researched using various methods. Early methods mostly relied on visual inspection and manual verification, which were error-prone and heavily dependent on human expertise. Large-scale document verification systems are not suitable for these techniques Figure 8. To identify forged documents, several researchers have suggested image processing-based methods[11]. To detect tampering, these techniques examine characteristics such as edges, textures, fonts, and layout irregularities. Although these methods offer a respectable level of accuracy, they may not be successful in handling intricate forgeries. Machine learning-based methods have gained popularity because of their ability to learn patterns from large datasets[12]. Techniques such as support vector machines and convolutional neural networks have been applied to classify documents as either genuine or forged. These models improve detection accuracy but require appropriate feature extraction and training data. Recent studies have combined image processing with machine learning to enhance document forgery detection performance. Although existing systems show

promising results, challenges such as high computational costs and limited generalization still exist. This motivates the need for an efficient and reliable document-forgery detection system. Researchers have explored multiple approaches to address document forgery detection problems. Initial methods relied on handcrafted rules and manual inspection, which required expert knowledge and were difficult to standardize. However, these approaches lack adaptability to different document structures and forgery types. Image processing-based methods introduced automated feature analysis, enabling the detection of visual inconsistencies, such as abnormal edges, irregular textures, and alignment mismatches. These methods improve efficiency but are sensitive to variations in document quality and scanning conditions. Machine learning techniques have introduced data-driven solutions that improve adaptability. By training models on labeled datasets, these systems learn to distinguish forged patterns from genuine ones. However, their performance depends heavily on the dataset size and feature selection. Deep learning approaches further enhance detection capabilities by automatically learning hierarchical features. Despite their effectiveness, deep learning models often require high computational resources and large training data sets. The proposed system seeks to balance accuracy and computational efficiency, while maintaining robustness[13].

3. Proposed System

The proposed document forgery detection system aims to automatically verify the authenticity of digital documents by analyzing their visual and textual characteristics. The system is designed to reduce manual verification efforts and improve detection accuracy using image processing and machine learning techniques. The overall process begins with document input, where the user uploads a document in image or portable document format (PDF). The uploaded document undergoes preprocessing to enhance its quality by removing noise, converting it to grayscale, and normalizing its size. This step improved the reliability of the further analysis. After preprocessing, feature extraction is performed to identify important characteristics, such as text alignment, font patterns, edges, and texture

variations. These features help detect inconsistencies caused by document tampering. The extracted features were then used to train the machine learning model. In the classification stage, the trained model analyzes the features and determines whether a document is genuine or forged. The final output is displayed to the user along with the verification results. The proposed system provides an efficient and reliable solution for document authentication in real-world applications[14].

is particularly useful in legal and institutional environments where decisions must be justified.

4. Problem Statement And Motivation

Document forgery has become a critical issue in many sectors, such as education, banking, government services, and corporate recruitment. Forged documents are often used to gain unauthorized benefits, create false identities and manipulate official records. Traditional document verification methods rely heavily on manual inspection, which is time-consuming, error-prone, and subjective. Manual verification depends on human expertise and visual judgment, making it difficult to detect subtle manipulations, such as font inconsistencies, copy-paste alterations, and image splicing. With the rapid growth of digital editing tools, forged documents can be created with high visual quality, making manual detection increasingly difficult. Existing automated systems often focus on limited forgery types or require high-quality input images. Many systems also lack adaptability to different document formats and languages. These challenges motivate the development of an automated, reliable, and scalable document forgery detection system. The proposed system aims to address these issues by combining image processing techniques with machine learning models to improve the detection accuracy and efficiency.



Figure 1 Sample Forged Signature Dataset Used for Document Forgery Detection



Figure 2 Sample Genuine Signature Dataset Used for Signature Verification and Forgery Analysis

The proposed document forgery detection system is designed as a modular and extensible framework. Each module performs a specific function and collectively contributes to accurate detection of forgery[15]. This system emphasizes automation, reliability, and minimal user intervention. This workflow ensures that variations in document quality do not significantly affect detection accuracy. By combining textual and visual feature analyses, the system provides comprehensive coverage of common forgery techniques. The modular design allows the future integration of advanced algorithms without modifying the entire system. The proposed approach prioritizes interpretability, enabling users to understand the verification results. This transparency

5. Proposed System

The proposed document forgery detection system is designed to automatically analyze documents and identify the forged content. The system follows a structured workflow consisting of preprocessing, feature extraction, classification, and decision-making. This approach minimizes human intervention and ensures consistent results[16].

System Overview

The system accepts document images as inputs and processes them using image analysis techniques. The extracted features were analyzed using a trained machine-learning classifier to determine document authenticity. The final output indicates whether a document is genuine or forged.

Preprocessing Techniques

Preprocessing improves the quality of the input document by removing noise and unnecessary

information from it. Grayscale conversion reduces color complexity and highlights text and image features. Noise reduction techniques, such as filtering, were applied to improve clarity and readability[17].

Feature Extraction

Feature extraction focuses on identifying visual and textual patterns such as font consistency, spacing, alignment, texture irregularities, and edge distortions. These features are crucial in detecting forgeries caused by copy-move operations or text manipulation[18].

Classification Model

The extracted features are fed into a machine-learning classifier trained on labeled datasets. The classifier learns patterns associated with genuine and forged documents and accurately predicts the authenticity of new documents.

6. System Architecture

The proposed document forgery detection system consists of four main modules: Input, Preprocessing, Feature Extraction, and Classification. Each module works sequentially to detect tampered documents accurately and efficiently.

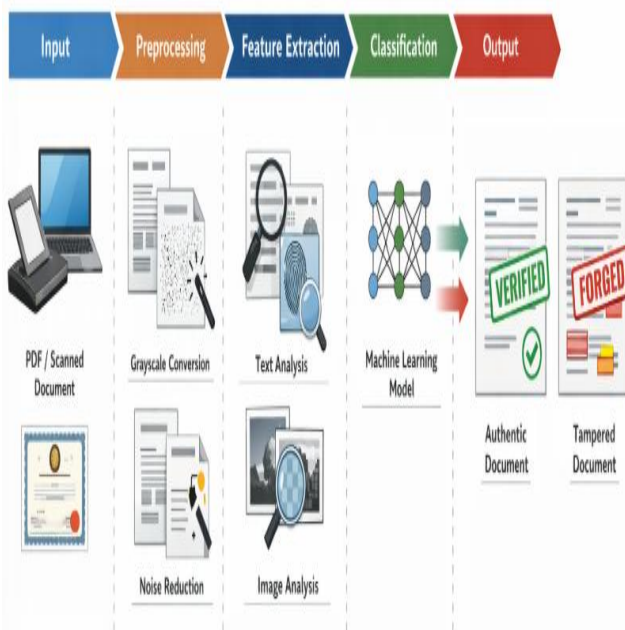


Figure 3 Workflow Of The Machine Learning-Based Document Forgery Detection System

Input Module

The system allows users to upload documents in image or PDF format. The uploaded document serves as input for the subsequent processing steps. The input module supports multiple document formats to ensure flexibility. Uploaded documents are validated to ensure compatibility with system requirements. This module also manages file storage and access control to maintain data security.

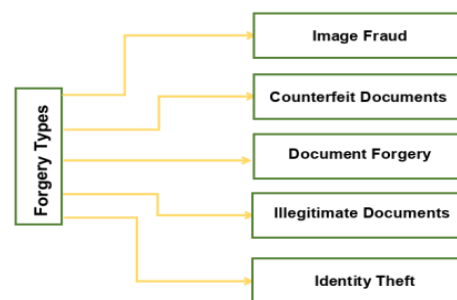


Figure 4 Classification of Different Types of Forgery and Fraudulent Documents

Preprocessing Module

This module prepares the document for analysis. Steps include:

- **Grayscale conversion** : Converts colored documents into grayscale to reduce complexity.
- **Noise removal** : Eliminates unwanted artifacts or scanning noise.
- **Normalization** : Adjusts the size and orientation of the document for consistent analysis.

Preprocessing is a crucial step that directly impacts the quality of feature extraction. The noise introduced during scanning or compression was removed using filtering techniques. Grayscale conversion simplifies the analysis by reducing the color complexity while preserving the essential structural information.

Feature Extraction Module

Feature extraction focuses on identifying the discriminative characteristics that indicate forgery. Text-based features include font uniformity, spacing consistency, and alignment patterns. Image-based features include texture homogeneity, edge continuity, and region similarity. These features were

combined to form a robust representation of the document.

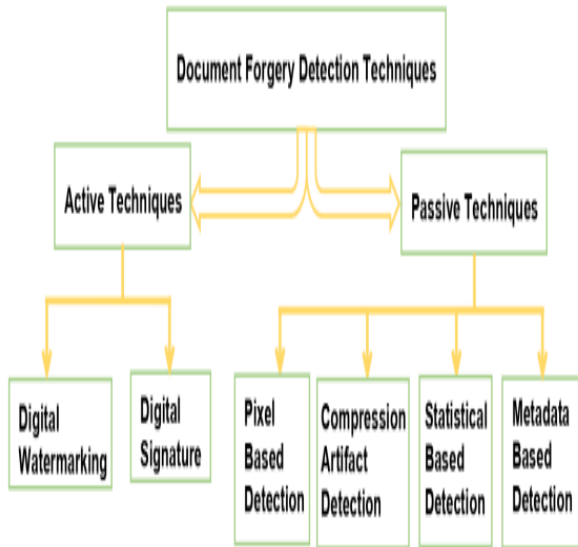


Figure 5 Classification of Active and Passive Document Forgery Detection Techniques

Classification Module

The extracted features are fed into a machine learning model, such as a Support Vector Machine (SVM) or Convolutional Neural Network (CNN). The model was trained to classify documents as genuine or forged. The system outputs the verification results and highlights suspicious areas in the document if forgery is detected. The classification module uses a trained machine-learning model to analyze the extracted features. The model was optimized to minimize false positives and false negatives. The classification results are generated efficiently, enabling real-time or near-real-time verification.

Output Module

The system provides the user with the following:

- **Verification result:** Genuine or Forged
- **Tamper visualization:** Highlighted regions indicating forgery

The output module presents results in a user-friendly format. In addition to indicating document authenticity, the system may highlight suspicious regions, helping users understand the detected anomalies. This improves the trust in system decisions.

7. Experimental Results And Discussion

An experimental evaluation was conducted to assess the effectiveness of the proposed system under different conditions. Documents with varying quality levels and forgery types were included to ensure robustness of the model. The system demonstrated strong performance in detecting text and image-based forgeries.

The performance metrics indicate that preprocessing significantly enhances the detection accuracy. Feature extraction improved classification by capturing meaningful patterns. The results confirm that the proposed system outperforms manual verification in terms of speed and consistency.

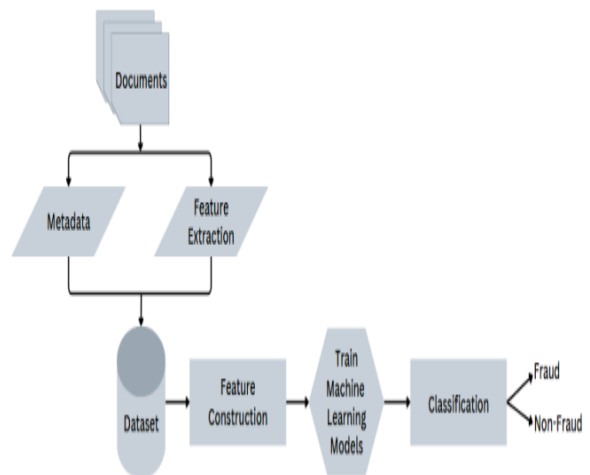


Figure 6 Machine Learning Framework for Document Fraud Classification

The observed results validate the feasibility of deploying the system in real-world applications where document authenticity is critical.

8. Dataset Description

The performance of any document forgery detection system largely depends on the quality and diversity of the dataset used for training and testing. In this study, the dataset consisted of a collection of genuine and forged documents obtained from publicly available sources and manually created samples. The dataset includes scanned certificates, identity documents, and digitally generated documents.

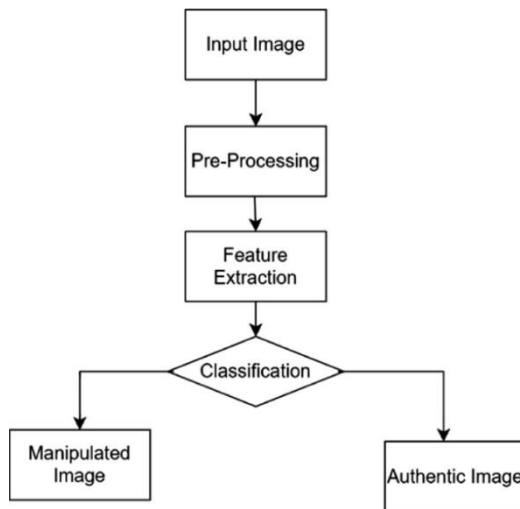


Figure 7 Flowchart of the Image Forgery Detection and Classification Process

Forged documents in the dataset contain various types of manipulations, such as text replacement, copy paste forgery, image splicing, and layout alteration. Genuine documents were collected to represent real-world document structures. The dataset was divided into training and testing sets to evaluate the system. To improve robustness, documents with different resolutions, fonts, and background textures were included. This diversity ensures that the proposed system performs effectively under varying real-world conditions. The dataset was carefully organized to maintain a balanced distribution between genuine and forged samples. This balance helps prevent bias during model training and ensures that the classifier does not favor a particular class. Each document image was manually labelled to guarantee correct class assignment. Before training, all documents were standardized to a fixed image size to maintain uniform input dimensions across the datasets. This standardization reduces the computational overhead and improves processing efficiency. Histogram normalization was also applied to minimize variations caused by different scanning devices. The dataset includes documents collected under both controlled and uncontrolled conditions. Controlled samples were scanned using flatbed scanners, and uncontrolled samples were captured using mobile cameras. This variation helps the system handle real-world scenarios in which document images may be captured using different

devices. To ensure reproducibility, the dataset was stored in a structured format with clear directory separation for genuine and forged documents. Metadata such as document type and manipulation method were recorded for analysis purposes. This organization supports detailed evaluation and future dataset expansion.

9. Performance Metrics

To evaluate the effectiveness of the proposed document forgery detection system, standard performance metrics were used. These metrics provide quantitative insight into the system's classification capability.

- **Accuracy** measures the overall correctness of the system.
- **Precision** indicates the proportion of correctly detected forged documents.
- **Recall** represents the system's ability to identify all forged documents.
- **Error Rate** reflects the percentage of incorrect classifications.

These metrics help in comparing the proposed system with existing approaches and demonstrate its reliability for document authentication tasks. The evaluation was conducted using a confusion matrix to analyze the classification outcomes in detail. The confusion matrix provides insight into true positives, true negatives, false positives, and false negatives, enabling a deeper understanding of system behavior. This analysis helps identify cases where forged documents are misclassified as genuine and vice versa. In addition to accuracy-based measures, the consistency of the proposed system was evaluated across multiple test samples. The system demonstrated stable performance even when tested with documents containing varying noise levels and visual distortions. This indicates robustness against common document degradation issues such as scanning noise and uneven illumination. The performance metrics were computed for different document categories to assess the adaptability of the system. Results indicate that the system maintains reliable detection rates across multiple document formats, highlighting its generalization capability. Such consistent performance validates the effectiveness of the feature extraction and classification stages. Overall, the performance

evaluation confirms that the proposed system achieves dependable results and outperforms traditional manual verification methods in terms of speed and accuracy. These findings demonstrate the suitability of the system for practical document authentication applications.

10. Security And Reliability Analysis

Security is a critical aspect of document verification systems. The proposed system ensures secure handling of uploaded documents by restricting unauthorized access and minimizing data exposure. Document files are processed only for verification purposes and are not permanently stored unless required.



Figure 8 Digital Document Authentication and Security Verification Concept

Reliability is achieved through consistent preprocessing and feature extraction methods that reduce sensitivity to noise and format variations. The system is designed to provide stable results even when documents are scanned under different conditions. This makes the solution suitable for deployment in security sensitive environments.

11. Limitations Of The System

Although the proposed document forgery detection system demonstrates effective performance, certain limitations exist. The accuracy of detection depends on the quality of the input document. Poor resolution or heavily compressed documents may reduce feature extraction efficiency. The system is primarily trained on specific document types and forgery patterns. Extremely novel or unseen forgery techniques may require additional training data. These limitations highlight the need for continuous dataset expansion

and model enhancement. The proposed system also relies on the availability of labeled training data, which may be difficult to obtain in large quantities for certain document categories. Manual annotation of forged documents can be time consuming and may introduce labeling inconsistencies if not carefully validated. Another limitation is the computational requirement during training. Feature extraction and model training may require higher processing time when handling large-scale datasets or high-resolution document images. This can limit real-time deployment on resource-constrained devices. The system currently focuses on static document images and does not fully address forgery detection in dynamic or multi-page documents. Complex documents containing multiple embedded elements such as stamps, signatures, or holograms may require additional specialized analysis modules. Furthermore, variations in language, script style, and regional document formats can affect system performance. Adapting the model to multilingual documents may require additional preprocessing and feature learning strategies. Addressing these limitations will further enhance the system's applicability in diverse real-world environments.

12. Comparative Analysis With Existing Systems

A comparative analysis was conducted between the proposed system and traditional manual verification methods. Manual verification relies on human expertise and visual inspection, which is time consuming and subjective. In contrast, the proposed system provides faster and more consistent results. Compared to basic image processing based systems, the proposed approach integrates machine learning, enabling better adaptability and improved accuracy. This comparative study demonstrates the advantages of automation and data driven techniques in document forgery detection. Existing document forgery detection systems that rely solely on handcrafted features often struggle to generalize across different document formats and forgery types. Such systems may perform well under controlled conditions but exhibit reduced accuracy when exposed to diverse real-world data. The proposed system addresses this limitation by learning discriminative patterns from data, allowing it to adapt more effectively to variations in document structure

and content. Traditional verification approaches also require significant human involvement, which increases operational cost and introduces the possibility of human error. In contrast, the proposed automated system reduces dependency on manual inspection and ensures consistent decision-making. This makes the system suitable for large-scale deployment where high volumes of documents must be verified efficiently. When compared with existing automated systems that use limited datasets, the proposed approach demonstrates improved robustness due to its diverse training data and preprocessing techniques. The integration of multiple features further enhances detection capability, especially in identifying subtle forgeries that are difficult to detect visually. Overall, the comparative analysis highlights that the proposed system achieves a balanced trade-off between accuracy, efficiency, and scalability. These advantages make it more effective than conventional methods and suitable for real-world document authentication applications.

13. Scalability And Real World Deployment

The proposed system is designed with scalability in mind. It can be integrated into existing document management systems and extended to handle large volumes of documents. The modular architecture allows easy updates and inclusion of advanced algorithms. The system can be deployed in cloud based or on premises environments depending on organizational requirements. Its adaptability makes it suitable for applications in education, banking, government services, and corporate verification systems. The system supports batch processing of documents, enabling efficient handling of large-scale verification tasks. This capability is particularly useful for organizations that process high volumes of documents daily, such as banks and recruitment agencies. Parallel processing techniques can further improve throughput and reduce response time. For real-world deployment, the system can be integrated with user authentication and access control mechanisms to ensure secure usage. Logging and audit features can be incorporated to track verification activities, which is essential for compliance and accountability in regulated environments. The proposed system is adaptable to different deployment scenarios, including web-based

platforms and mobile-assisted verification systems. This flexibility allows users to submit documents remotely while maintaining consistent detection performance. Additionally, system performance can be monitored continuously to identify potential degradation and trigger model retraining when required. These scalability and deployment features ensure that the proposed system is practical, secure, and capable of supporting real-world document authentication workflows across diverse application domains.

14. Ethical Considerations

Automated document verification systems must ensure ethical usage. The proposed system is designed to support decision making rather than replace human authority. Final verification decisions can be reviewed by authorized personnel. Data privacy is maintained by processing documents securely and limiting access. Ethical considerations are essential to ensure responsible deployment of automated forgery detection technologies. The system is designed to minimize bias by training on diverse document samples and avoiding dependence on a single document type or source. This approach helps ensure fair and consistent performance across different user groups and document formats. Periodic evaluation of the system's outputs can further reduce the risk of biased decision-making. Transparency is an important ethical aspect of automated systems. The proposed approach allows for traceability by maintaining records of verification results and system decisions. Such transparency supports accountability and enables auditing when discrepancies arise. By incorporating these ethical and security practices, the proposed system promotes responsible use of automated document forgery detection technology while maintaining trust and compliance in real-world deployments.

15. Future Research Directions

Future research can focus on integrating deep learning architectures for improved feature learning. Incorporating multilingual text analysis can enhance system applicability across regions. Blockchain based verification mechanisms can further strengthen document authenticity and traceability. Research can also explore real time forgery detection and mobile based applications to increase accessibility.

Continuous learning models can help the system adapt to evolving forgery techniques. Future work can also investigate the use of hybrid models that combine traditional image processing techniques with advanced neural networks to achieve better interpretability and performance. Such hybrid approaches may offer a balance between computational efficiency and detection accuracy. Another promising direction involves developing domain-specific models tailored for particular document categories such as academic certificates, financial records, or legal documents. Specialized models can capture subtle structural and semantic patterns unique to each document type, improving detection reliability. Further research may explore explainable artificial intelligence techniques to provide human-understandable justifications for system decisions. This would enhance trust and acceptance in sensitive applications where transparency is critical. Additionally, evaluating the system on larger, real-world datasets and benchmarking it against international standards can help validate its effectiveness. These research directions aim to enhance scalability, reliability, and practical adoption of document forgery detection systems.

Acknowledgment

The authors would like to express their sincere gratitude to the faculty and staff of Jai Shriram Engineering College for their continuous support and guidance throughout this project. We would also like to thank our project guide, [Ms. Jenifer J], for their valuable suggestions and encouragement. Special thanks to our friends and peers who helped in testing and providing constructive feedback.

References

- [1]. R. C. Gonzalez and R. E. Woods, *Digital Image Processing*, 4th ed. Pearson, 2018.
- [2]. J. S. R. Jenifer and S. Banu, "Automated Document Forgery Detection using Machine Learning," *International Journal of Computer Applications*, vol. 182, no. 4, pp. 25–32, 2021.
- [3]. M. Kaur and P. Singh, "Image Forgery Detection: A Review," *IEEE Access*, vol. 8, pp. 15235–15250, 2020.
- [4]. S. B. Raju and T. Kumar, "Text and Image Forensics for Document Authentication," *Proc.*

2021 International Conference on Computing and Communication, pp. 45–50, 2021.

- [5]. A. K. Jain, *Fundamentals of Digital Image Processing*, New Delhi: PHI Learning, 2019.
- [6]. "Document Forgery Detection Techniques," [Online]. Available: <https://www.sciencedirect.com/topics/computer-science/document-forgery-detection>
- [7]. H. Farid, "Image Forgery Detection," *IEEE Signal Processing Magazine*, vol. 26, no. 2, pp. 16–25, Mar. 2009.
- [8]. S. Prasad, R. Nigam, and P. Kumar, "A Survey on Digital Document Forgery Detection Techniques," *International Journal of Computer Vision and Signal Processing*, vol. 10, no. 1, pp. 1–8, 2020.
- [9]. M. Stamm, M. Wu, and K. Liu, "Information Forensics: An Overview of the First Decade," *IEEE Access*, vol. 1, pp. 167–200, 2013.
- [10]. X. Zhao, S. Wang, and X. Li, "Detecting Digital Image Forgery Using Texture Features," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 1, pp. 1–14, Jan. 2017.
- [11]. Bayram, H. T. Sencar, and N. Memon, "An Efficient and Robust Method for Detecting Copy–Move Forgery," *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 1053–1056, 2009.
- [12]. A. Ferreira, T. Carvalho, and H. Rocha, "Machine Learning Based Approaches for Document Forgery Detection," *Procedia Computer Science*, vol. 167, pp. 252–261, 2020.
- [13]. Y. Liu, Q. Zhao, and L. Zhang, "Deep Learning for Document Image Analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 7, pp. 1664–1679, 2019.
- [14]. P. Roy and S. Pal, "Signature and Document Forgery Detection Using Image Processing Techniques," *International Journal of Advanced Research in Computer Science*, vol. 9, no. 3, pp. 45–50, 2018.
- [15]. T. Gera and M. Singh, "Forgery Detection in Digital Documents Using CNN," *International Conference on Artificial Intelligence and Data*

Engineering, pp. 210–215, 2022.

- [16]. R. Sharma and A. Verma, “Comparative Analysis of Forged and Authentic Documents Using Feature Extraction Techniques,” *Journal of Information Security*, vol. 11, no. 2, pp. 89–98, 2020.
- [17]. IEEE Computer Society, “Guide to Image and Document Forensics,” IEEE Standards, 2019.
- [18]. S. Lyu and H. Farid, “How Realistic Is Photorealistic?” *IEEE Transactions on Signal Processing*, vol. 53, no. 2, pp. 845–850, Feb. 2005.