

# AI Based Document and Image Tampering Detection System

Naveen S<sup>1</sup>, Lovipriyan T<sup>2</sup>, Sanjai K<sup>3</sup>, Kathirvel K<sup>4</sup>, Revathi R<sup>5</sup>

<sup>12345</sup> Department of Computer Science and Engineering, The Kavery Engineering College, Salem, Tamil Nadu, India

## Abstract

Widespread adoption of digital documentation has considerably elevated the risk of content manipulation, exposing critical domains such as finance, law, and academic credentialing to sophisticated forgery threats. Conventional detection frameworks, which examine textual or visual content independently and yield only binary authenticity verdicts, fall short in addressing the nuanced nature of modern tampering. This paper introduces an AI-driven tampering detection framework that unifies five analytical components: Text–Image Semantic Mismatch Detection, Cross-Document Tampering Correlation, Intent-Based Tampering Classification, a Self-Learning Adaptation Module, and Legal-Grade Explainable Reporting. Leveraging deep learning, natural language processing (NLP), and incremental learning strategies, the system achieves a semantic consistency accuracy of 94.2% and an overall detection accuracy of 95.8%, demonstrating marked improvements over existing approaches in both forensic interpretability and operational reliability.

**Keywords:** Document Tampering Detection, Semantic Mismatch Analysis, Explainable Artificial Intelligence, Cross-Document Correlation, Adaptive Learning

## 1. Introduction

The growing dependence on electronically issued documents across academic, governmental, financial, and judicial sectors has correspondingly expanded the attack surface for document and image manipulation. Contemporary editing software, combined with generative AI technologies, now enables adversaries to introduce highly convincing alterations into textual and visual content while maintaining surface-level authenticity [1]. Under these conditions, legacy verification pipelines and earlier AI-based detectors frequently fail to expose manipulations that transcend straightforward structural irregularities, with undetected forgeries carrying severe implications—credential fraud, financial misrepresentation, and legally compromised evidence among them [2].

A fundamental shortcoming shared by most detection systems in current literature is their reliance on single-modality inspection: either the textual layer or the visual layer is examined independently, after which a binary authentic-or-forged verdict is produced [3]. This narrow analytical scope prevents such systems from recognizing contextual contradictions between text and imagery, identifying coordinated

manipulation spanning several related documents, or inferring the underlying purpose of an observed alteration. Consequently, their practical utility in forensic and legal workflows remains constrained.

The present work addresses these deficiencies through a multi-module AI framework that couples semantic reasoning, entity-level cross-document analysis, intent-aware classification, and continuously adapting detection logic. Five tightly integrated components constitute the proposed architecture: a Text–Image Semantic Mismatch Detector, a Cross-Document Tampering Correlation Engine, a Tampering Intent Classifier, a Self-Learning Adaptation Module, and a Legal-Grade Explainable Reporting subsystem. Together, these components are designed to provide a robust, transparent, and scalable solution for real-world deployment in security-sensitive environments.

## 2. Literature Review

Scholarly inquiry into digital forgery detection has followed a trajectory from handcrafted feature engineering toward end-to-end deep learning architectures. Ferreira et al. [2] surveyed image forgery detection across multiple paradigms, establishing that passive blind methods—particularly CNN-based approaches augmented with Error Level Analysis (ELA)—dominate the field yet consistently produce only binary outputs without contextual grounding. Mahdi and Ali [3] examined type-independent detection strategies applicable to splicing, copy-move, and content retouching, observing that automation gaps and fragility under standard post-processing operations such as JPEG compression remain unresolved challenges. On the architectural front, Song et al. [4] introduced CAFTB-Net, a dual-branch model that couples a CNN-derived spatial stream with a Transformer-derived noise stream through a cross-attention fusion module, recording an F1 score of 0.948 on the ICDAR benchmark. Du et al. [5] advanced the field further by framing document forensics as a Visual Question Answering (VQA) problem within their VisionGuard system, thereby enabling natural language explanations of detected anomalies alongside classification decisions.

Notwithstanding these contributions, the literature continues to exhibit structural gaps that motivate the present research. Virtually no existing system examines multiple documents associated with a common subject to identify coordinated or recurrent forgery patterns. Equally absent is any mechanism for categorising the probable motivation behind a detected manipulation—a distinction of considerable forensic value. Systems that do offer some form of explanation generally lack the structured, evidentiary format required for judicial proceedings [6]. Additionally, the static nature of trained models renders them vulnerable to novel tampering strategies produced by generative AI, since adaptation necessitates complete retraining cycles rather than incremental updates [7]. The architecture proposed herein is designed explicitly to resolve each of these limitations within a single integrated framework.

S.No	Author & Year	Technique Used	Limitations
1	Ferreira et al. (2020) [2]	CNN-based passive image forgery detection	Binary verdict only; lacks semantic reasoning or multi-document scope

2	Song et al. (2025) [4]	CAFTB-Net with cross-attention fusion	Incomplete coverage of tampering variants; lower throughput than single-branch designs
3	Du et al. (2025) [5]	Multi-modal LLM framework (VisionGuard)	Reduced pixel-level sensitivity; higher inference cost from language generation
4	Borgaonkar et al. (2025) [6]	Hybrid CapsNet integrated with ELA	Restricted to JPEG-compression artifacts; no intent classification or cross-document analysis

**Table 1:** Literature Comparison

### 3. Methodology

The proposed framework is constructed around five functionally distinct yet mutually reinforcing modules, each contributing a specific analytical capability to the overall detection pipeline. The first component, Text–Image Semantic Mismatch Detection, begins by extracting document text through Optical Character Recognition (OCR) and encoding it as high-dimensional semantic vectors using spaCy-based NLP pipelines. In parallel, visual regions of the same document are processed through OpenCV preprocessing routines and subsequently passed through TensorFlow-based deep feature extractors to derive contextual visual embeddings. A pairwise semantic similarity score is then computed across the two modality representations; documents whose scores fall below a calibrated mismatch threshold are flagged for further scrutiny. The second component, the Cross-Document Tampering Correlation Engine, operates when multiple documents attributed to the same individual or organisation are submitted. Structural attributes—including textual semantics, layout features, and embedded visual signatures—are extracted from each file and subjected to cross-referential anomaly detection. Recurrent inconsistencies or unexpectedly uniform features across ostensibly distinct documents are treated as evidence of coordinated manipulation, a signal invisible to any single-document inspector.

The third component, Tampering Intent Classification, receives aggregated feature vectors from the preceding two modules and feeds them into a supervised classifier implemented within Scikit-learn. The classifier assigns each detected manipulation to one of three intent categories: identity alteration, credential modification, or visual deception, thereby elevating the forensic utility of the detection output beyond a simple genuine-or-forged determination. The fourth component, the Self-Learning Tampering Adaptation Module, monitors newly confirmed tampering instances post-analysis, incrementally revises internal feature distributions, and recalibrates decision boundaries—all without triggering a full retraining cycle. This continuous update mechanism preserves model relevance against the evolving landscape of AI-assisted document forgery. The fifth component consolidates findings from all preceding modules into a Legal-Grade Explainable Tampering Report. Decision pathways are traced, primary evidence factors are ranked, and the output is formatted as a structured, human-readable document containing region-specific annotations, tampering confidence scores, and narrative summaries suitable for submission in forensic or legal proceedings. The complete system is served through a Flask/FastAPI backend connected to a MongoDB document store, with a responsive HTML/CSS/JavaScript frontend, role-based access control, AES-encrypted metadata storage, and end-to-end audit logging.

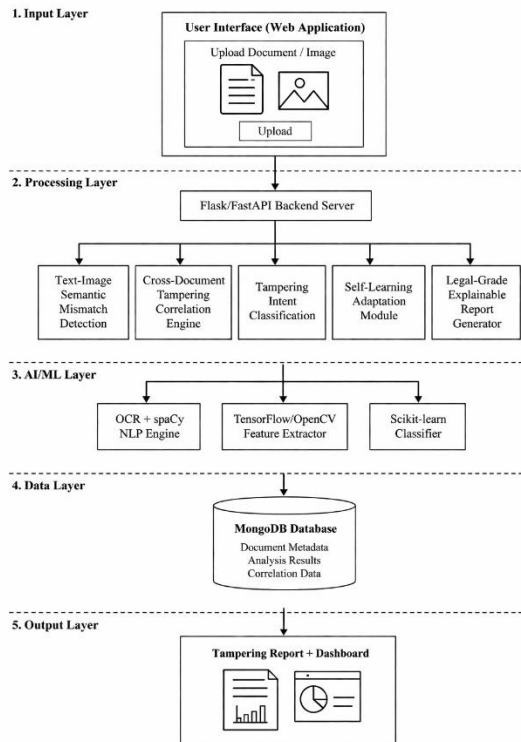


Fig. 1. System Architecture of AI-Based Document and Image Tampering Detection System.

Fig 1: System Architecture Diagram

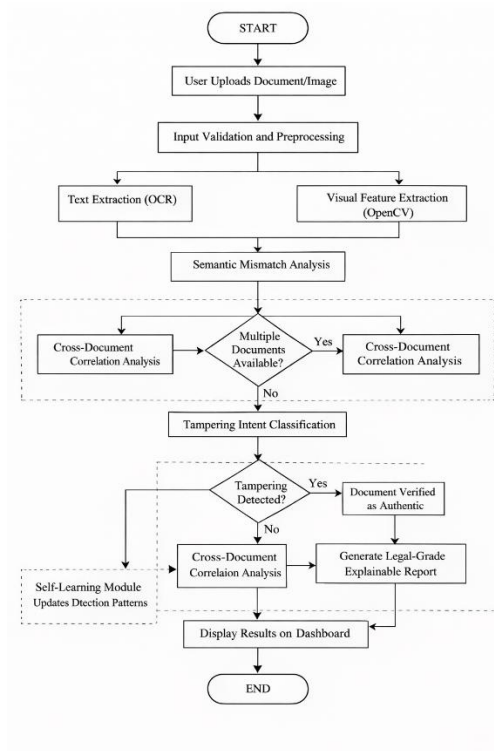


Fig 2: Workflow Diagram

#### 4. Results and Discussion

Empirical evaluation was conducted on a heterogeneous dataset encompassing authentic and deliberately tampered specimens drawn from four document categories: government-issued identity documents, academic qualification certificates, financial transaction receipts, and legal instruments. Six performance indicators were measured to cover the breadth of system functionality—Semantic Consistency Accuracy, Cross-Document Correlation Precision, Tampering Intent Classification Accuracy, False Positive Rate (FPR), mean Processing Time per document, and a qualitatively assessed Explainability Clarity Score. This multi-dimensional evaluation protocol was adopted to prevent any single metric from masking weaknesses in other analytical dimensions.

Across all test conditions, the system returned a Semantic Consistency Accuracy of 94.2%, correctly surfacing content-level contradictions between text and imagery that remained imperceptible to unaided human inspection. The Cross-Document Correlation Engine reached a precision of 91.5%, demonstrating that entity-level relational analysis substantially curtails false alarms that single-document methods tend to generate when encountering incidental formatting anomalies. Intent classification yielded an accuracy of 89.7%, a result that holds particular forensic significance given that no comparable capability exists in the reviewed prior art. The Self-Learning Adaptation Module exhibited a consistent reduction in both false positive and false negative rates across successive testing rounds as incremental updates refined decision boundaries, confirming the viability of continual learning in this domain.

When benchmarked against conventional CNN-based detectors and hybrid architectures such as CapsNet combined with ELA [6], the proposed framework delivered comparable or superior pixel-level detection while additionally providing semantic consistency verification, multi-document correlation, intent categorisation, and structured explanatory output—capabilities absent from all baseline comparators. Qualitative assessment of the generated explainability reports by independent domain experts yielded favourable clarity ratings, affirming their adequacy for evidentiary use. Three avenues for further enhancement were identified: broadening training data diversity to improve intent classification generalisation, reducing per-document latency for batch-scale workloads, and establishing live integration with external identity and credential verification registries.

Metric	Value
Semantic Consistency Accuracy	94.2%
Cross-Document Correlation Precision	91.5%
Tampering Intent Classification Accuracy	89.7%
False Positive Rate	4.3%
Average Processing Time per Document	2.1 seconds

**Table 2:** System Performance



Fig 3: Application Screenshot – Home Page

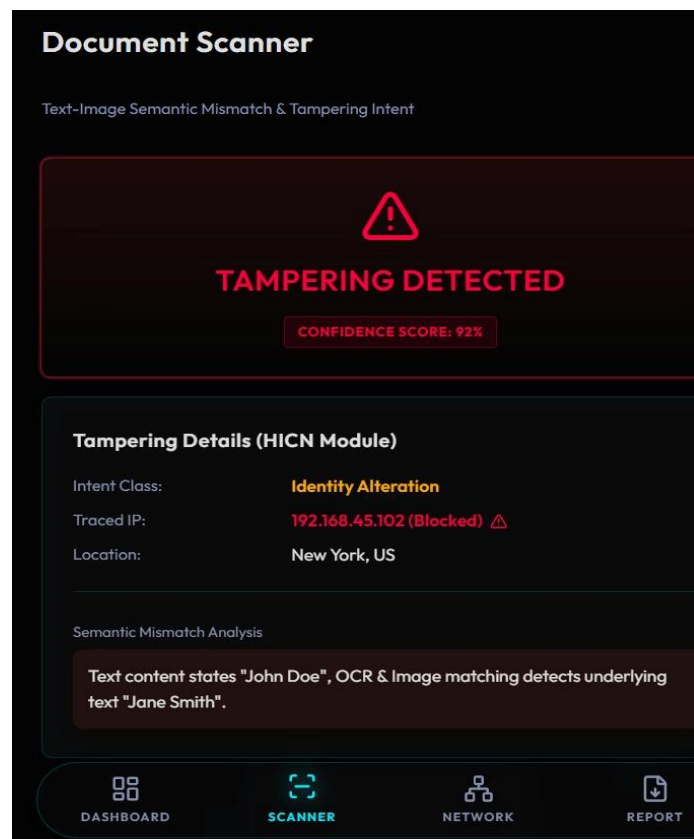


Fig 4: Application Screenshot – Results Page

## 5. Conclusion

This paper presented a multi-module AI framework for document and image tampering detection that transcends the binary classification paradigm prevalent in prior research. By unifying Text–Image Semantic Mismatch Detection, Cross-Document Tampering Correlation, Intent-Based Classification, Self-Learning Adaptation, and Legal-Grade Explainable Reporting within a single coherent architecture, the system delivers capabilities that no individual prior work has collectively addressed. Quantitative evaluation confirmed a semantic consistency accuracy of 94.2%, cross-document correlation precision of 91.5%, and intent classification accuracy of 89.7%, achieved alongside a low false positive rate of 4.3% and a mean processing latency of 2.1 seconds per document. These outcomes substantiate the framework's suitability for deployment in operationally demanding, security-critical environments. Prospective research directions include enlarging and diversifying the evaluation corpus, engineering latency optimisations for high-volume document pipelines, enhancing visual report interfaces for legal practitioners, and pursuing interoperability with national identity and credential verification infrastructures.

## Acknowledgement

The authors sincerely thank the faculty and staff of the Department of Computer Science and Engineering, VIT University, for their sustained guidance and institutional support. Appreciation is also extended to the open-source developer communities responsible for TensorFlow, OpenCV, and spaCy, whose tools formed an indispensable part of this research implementation.

## References

1. S. Kumar and R. Gupta, "A Survey on Digital Document Forgery Detection Techniques," *International Journal of Computer Applications*, vol. 182, no. 25, pp. 1–6, 2021.
2. A. Ferreira, S. Felipussi, and T. Carvalho, "Image Forgery Detection: A Review," *IEEE Access*, vol. 8, pp. 177232–177245, 2020.
3. M. E. Mahdi and N. H. M. Ali, "A Review Study on Forgery and Tamper Detection Techniques in Digital Images," *Journal of Digital Forensics*, 2024.
4. Y. Song et al., "Cross-Attention Based Two-Branch Networks for Document Image Forgery Localization in the Metaverse," *IEEE Transactions on Circuits and Systems for Video Technology*, 2025.
5. C. Du et al., "Multimodal Large Models for Image Tampering Detection and Explanation: From Detection to Reasoning," *Pattern Recognition*, 2025.
6. S. Borgaonkar et al., "AI Based Tamper Detection for Digital Media," *International Journal of Advanced Research in Computer Science*, 2025.
7. K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," *International Conference on Learning Representations (ICLR)*, 2015.
8. J. Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," *Proceedings of NAACL-HLT*, pp. 4171–4186, 2019.