

# AI-Powered Fake Certificate and QR Code Verification System

VISHWANATHAM HARSHITHA<sup>1</sup>, SAMALA ASHRITHA<sup>2</sup>, MIDIVELLI ARAVINDA SWAMI<sup>3</sup>,  
UDUTHA AKSHAY SAI<sup>4</sup>, REGULAPATI AKHILA RAO<sup>5</sup>

<sup>1,2,3,4</sup>Department of Information Technology, J.B. Institute of Engineering & Technology, Hyderabad

<sup>5</sup>Assistant Professor, Department of Information Technology, J.B. Institute of Engineering & Technology, Hyderabad

**Abstract**—The rapid spread of digital credentials across academic institutions, recruitment organisations, and online service platforms has created a pressing demand for dependable certificate authentication tools. Traditional methods of manual verification are inherently slow, subject to human error, and increasingly ineffective against sophisticated document forgeries and tampered QR codes. This paper presents an AI-powered web-based system that addresses these shortcomings by automatically analysing uploaded certificate images and QR code images to determine their authenticity.

The proposed system employs Convolutional Neural Network (CNN) models to extract deep visual features from both certificate documents and QR code patterns. Two dedicated CNN pipelines run in parallel—one for certificate analysis and one for QR code validation—and their outputs are fused in a prediction module that produces a binary Real/Fake classification, a fake-probability confidence score ranging from 0 to 100, and a human-readable textual explanation of the detected anomalies. The backend is built with the Django framework, while the frontend leverages React with TypeScript, HTML, CSS, and JavaScript. Core machine learning operations rely on TensorFlow, Keras, OpenCV, NumPy, and scikit-learn.

Experimental results drawn from prior related work on comparable tasks consistently demonstrate CNN-based approaches achieving detection accuracies between 93% and 98.6%, validating the soundness of the chosen architecture. The system delivers verification outcomes in near real-time, supports downloadable and shareable reports, and offers an interpretable interface that makes the authentication process transparent to non-technical users. The proposed solution is intended to help universities, employers, and certification bodies reduce manual overhead, curb credential fraud, and build trust in digital qualification ecosystems.

**Index Terms**—Fake certificate detection, QR code tampering, convolutional neural network, deep learning, document forgery, image classification, fake-probability score, Django, web-based verification

## I. INTRODUCTION

Academic and professional credentials act as the primary bridge between an individual's achievements and the opportunities they pursue. Whether a student is applying for graduate admission, a job-seeker is submitting a portfolio, or a contractor is demonstrating compliance licences, the institution on the receiving end must be able to trust that the document in front of them is genuine. Yet this trust is under growing threat. Advances in image-editing tools, freely available desktop publishing software, and high-resolution printing have collectively made it far easier than ever to produce convincing

forgeries of degree certificates, professional qualifications, and training completion cards. Equally concerning is the rise of tampered QR codes embedded within documents—codes that superficially scan as valid but either redirect to fraudulent verification pages or carry falsified payload data [5].

Institutions typically respond to this threat through one of three routes: manual visual inspection, rule-based OCR systems, or dedicated cryptographic databases. Manual inspection is inherently unreliable at scale; a trained officer may catch obvious anomalies in layout or typography, but will miss subtle pixel-level manipulations that modern editing tools leave behind. Rule-based OCR approaches, as studied by Patel et al. [9], can automate the extraction of key fields such as student name, roll number, and grade, but they achieve only moderate accuracy because they fail on noisy, low-contrast, or non-standard certificate templates. Cryptographic hash solutions, explored by Boonkrong et al. [1], provide strong integrity guarantees when both the issuer and the verifier share a common infrastructure, but they break down entirely when a legacy certificate predates the hash registry or when an institution has not yet joined the scheme.

Deep learning approaches, particularly CNN-based architectures, have demonstrated compelling results across image-forensics tasks [4], [8], [11]. A CNN trained on authentic and forged certificate images learns to detect subtle irregularities in layout geometry, font consistency, seal placement, and compression artefacts that are invisible to the human eye but statistically detectable at the feature-map level. Extending this idea to QR code authenticity, researchers have shown that texture-embedding methods [2] and dual-branch multi-scale feature fusion networks [3] can discriminate between original and copied QR patterns with accuracies exceeding 95%.

Despite this body of work, a practical gap remains. Existing systems tend to address either certificate forgery or QR tampering in isolation. Few combine both verification pathways into a single, unified web interface that a non-expert user can operate without any specialised training. Furthermore, the majority of published prototypes stop at a binary decision; they do not generate a quantitative confidence score or offer a natural-language explanation of why a particular document was flagged—features that are essential when a human reviewer must act on the system's output.

This paper proposes a system, hereafter referred to as the

*AI-Powered Fake Certificate and QR Code Verification System*, that fills these gaps. The key contributions of this work are as follows.

- 1) A dual-CNN architecture that independently processes certificate images and QR code images, then routes predictions through a unified prediction module.
- 2) A probabilistic output layer that converts raw model logits into a fake-probability score in the range [0, 100], giving reviewers a graded indication of forgery risk rather than a hard binary label.
- 3) An explainability layer that maps the network's confidence back to natural-language descriptions of detected anomalies, making the verdict interpretable for administrative staff without a machine learning background.
- 4) A full-stack web deployment using Django on the backend and React + TypeScript on the frontend, enabling immediate practical adoption without any client-side software installation.

The remainder of this paper is structured as follows. Section II surveys the related literature. Section III analyses the limitations of existing verification systems. Section IV presents the proposed methodology and mathematical foundations. Section V covers implementation details. Section VI presents results and discussion. Section VII concludes the paper and outlines future directions.

## II. LITERATURE SURVEY

The challenge of verifying the authenticity of academic and professional documents has attracted considerable attention from the research community, leading to a diverse set of approaches spanning classical machine learning, deep learning, cryptographic hashing, and optical character recognition. This section reviews the most directly relevant prior work, grouping studies by their primary technique, and identifies the limitations that motivated the present research.

Kaur and Singh [7] investigated the use of Error Level Analysis (ELA) combined with classical machine learning classifiers for detecting tampered images. ELA works by recompressing a JPEG image at a known quality level and measuring the per-pixel difference between the original and the recompressed version; regions that were edited after the initial save exhibit noticeably different error levels. While the authors reported around 90% accuracy on their test set, the technique is inherently fragile: a skilled forger can neutralise ELA signatures by saving the document multiple times or by using PNG rather than JPEG as the final format. The approach also provides no output confidence measure, making it difficult for a human reviewer to calibrate trust in the result.

Qian et al. [6] moved beyond pixel-level analysis by proposing a hybrid model that uses a CNN for hierarchical feature extraction and an SVM as the final classifier. Trained on certificate images, this pipeline achieved a reported accuracy of approximately 98.6%, demonstrating that learned deep features capture layout regularities and typographic consistency far more robustly than hand-crafted descriptors. However, the model considered only the visual surface of the document; it

did not verify textual fields, embedded metadata, or security elements such as watermarks and digital signatures, leaving a clear attack surface for adversaries who preserve the visual appearance while altering specific data fields.

Chen et al. [4] addressed multi-region manipulation in scanned documents using a CNN-based splicing detector that achieves around 94% accuracy. The network learns to identify boundaries between authentic and inserted regions by analysing statistical inconsistencies in compression noise, lighting gradients, and edge profiles. While this approach generalises better than ELA to diverse forgery types, it struggles when an attacker edits multiple non-contiguous regions that individually appear clean, because the model lacks a global reasoning mechanism that can reconcile evidence across the entire page.

Boonkroong et al. [1] proposed a system that combines cryptographic hash verification with a machine learning layer for academic document authentication. Each issued certificate is registered in a central database alongside its cryptographic fingerprint; at verification time, the submitted document's fingerprint is recomputed and compared against the registry. The machine learning module then assesses documents whose fingerprints are absent from the registry, attempting to infer authenticity from structural features alone. The primary limitation is template dependency: the ML module was trained only on a fixed set of certificate layouts, and performance degraded significantly on certificates from institutions not represented in the training corpus.

In the QR code domain, Wang et al. [2] introduced a texture-hidden anti-counterfeiting QR code scheme in which subtle, imperceptible texture patterns are embedded into QR code images at generation time. At verification, a deep learning model attempts to recover and validate these hidden textures. The method achieved accuracy between 95% and 100% on the authors' proprietary dataset, but it requires that the original QR code be generated using the authors' specific encoding pipeline—meaning it cannot authenticate legacy QR codes already in circulation.

Zhang et al. [3] tackled copy-forgery attacks on anti-counterfeiting QR codes through DMF-Net, a dual-branch multi-scale feature fusion network. One branch extracts fine-grained local texture features, while the other captures global structural patterns; the branches are merged via an attention mechanism before the final classification head. DMF-Net achieved around 96% accuracy on their benchmark but showed sensitivity to lighting variation and photographic noise, suggesting that robustness to real-world capture conditions requires further study.

Harley et al. [11] evaluated several deep CNN architectures on the RVL-CDIP document image classification benchmark, establishing baseline performance figures for general document understanding. Although this work was not focused on forgery detection, it confirmed that CNNs pretrained on large document corpora learn transferable representations useful for downstream document analysis tasks, a finding that has since been widely exploited in the forgery-detection literature.

Hafemann et al. [10] addressed the related problem of offline handwritten signature verification using Siamese CNNs. A Siamese architecture learns an embedding space in which genuine signature pairs are pulled together and genuine-forgery pairs are pushed apart. The approach achieved high similarity accuracy but requires a database of known genuine signatures for each individual—a prerequisite that is not always satisfied in practice, particularly for legacy certificates where the original signatory's pen strokes were never digitised.

Sharma et al. [8] applied a straightforward CNN classifier to scanned document images for fraud detection, reporting 93% accuracy overall. The model's performance degraded on low-resolution scans, highlighting the importance of robust image preprocessing pipelines—a consideration that directly informs the preprocessing design in the current work.

Patel et al. [9] designed an OCR-based certificate verification system that extracts text fields using optical character recognition and then validates them against a rule-based knowledge base encoding expected patterns for dates, registration numbers, and institutional names. The system achieved only moderate accuracy because OCR quality is sensitive to scan resolution, font diversity, and document damage, and because rule bases require continual manual maintenance as certificate formats evolve.

Li et al. [5] demonstrated that CNN-based anomaly detectors can identify QR code tampering with approximately 95% accuracy by learning the statistical signature of authentic module patterns. The work noted, however, that the model does not inspect the decoded payload content, leaving it vulnerable to attacks that preserve the visual QR structure while modifying only the encoded URL or data string.

Taken together, these studies reveal three persistent gaps in the existing literature: (1) no published system jointly verifies both certificate images and QR codes within a single end-to-end pipeline; (2) quantitative confidence scores and natural-language explanations are rarely provided, reducing the practical utility of the systems for human reviewers; and (3) template dependency and lighting sensitivity remain open challenges for real-world deployment. The system proposed in this paper is designed to close all three gaps.

### III. EXISTING SYSTEM

Prior to the introduction of deep learning-based approaches, certificate and QR code verification relied on several conventional paradigms. This section categorises and critically examines these existing systems to motivate the need for the proposed architecture.

#### A. Manual Visual Inspection

The most widely deployed verification method remains manual examination by a trained administrative officer. The reviewer visually inspects the document for inconsistencies in typography, font weight, seal placement, border alignment, and overall layout quality. While experienced staff can detect crude forgeries, this approach has fundamental limitations: (a) it does not scale to high-volume verification environments such

as national recruitment boards or large universities; (b) it is susceptible to fatigue errors; and (c) it cannot detect pixel-level manipulations introduced by modern image-editing software, which are imperceptible to the naked eye.

#### B. OCR-Based Rule Systems

Systems of this category, exemplified by the work of Patel et al. [9], use Optical Character Recognition to extract structured text fields—student name, registration number, issue date, and grade—from the scanned certificate. The extracted values are then validated against a rule base that encodes expected formats and permissible value ranges. The key weaknesses of OCR-based systems are threefold. First, recognition accuracy degrades significantly on low-resolution scans, non-standard fonts, and physically damaged documents. Second, the rule base must be manually maintained and expanded every time a new issuing institution or certificate layout is introduced. Third, these systems have no mechanism to detect visual or structural forgeries that preserve all text content while altering graphical elements such as seals, watermarks, or signatories.

#### C. Error Level Analysis with Classical Machine Learning

Error Level Analysis (ELA) [7] detects tampered regions by exploiting the differential compression artefacts introduced when a JPEG image is re-saved after editing. Forged regions typically display higher error levels than the surrounding unmodified areas. Classical classifiers such as Support Vector Machines (SVM) or Random Forests are then trained on ELA feature maps to predict forgery. Although this technique achieves around 90% accuracy on controlled datasets, it is brittle in practice because an attacker can suppress ELA signatures simply by saving the final document as a lossless PNG file or by performing multiple JPEG re-compression passes. Furthermore, ELA provides no confidence score and offers no actionable explanation of where the manipulation occurred.

#### D. Cryptographic Hash-Based Systems

Systems such as that proposed by Boonkrong et al. [1] generate a cryptographic fingerprint of each certificate at issuance time and store it in a central registry. At verification time, the submitted document's fingerprint is recomputed and compared against the registry record. This approach provides strong integrity guarantees for documents registered in the system. However, it suffers from two critical dependencies: (a) the issuing institution must participate in the registry infrastructure, which excludes the vast majority of legacy certificates already in circulation; and (b) even a single-pixel alteration to the document image will invalidate the hash, causing the system to flag genuinely authentic documents that have undergone incidental file-format conversion or lossless compression.

#### E. Single-Domain CNN Classifiers

More recent approaches apply convolutional neural networks directly to certificate or QR code images. Qian et al. [6]

reported 98.6% accuracy using a CNN + SVM pipeline trained exclusively on certificate images. Chen et al. [4] achieved around 94% on a splicing detection task. Sharma et al. [8] obtained 93% on general document fraud detection. In the QR code domain, Li et al. [5] and Zhang et al. [3] independently demonstrated that domain-specific CNNs can identify QR tampering with accuracies of approximately 95% and 96%, respectively. While these results are encouraging, all of the above systems share three architectural shortcomings: (a) they address either certificate verification or QR code verification in isolation and cannot handle both within a unified pipeline; (b) they output a hard binary label without a quantitative confidence score, limiting their utility when a human reviewer needs to make a risk-calibrated decision; and (c) none provides a natural-language explanation of the detected anomalies, making the verdict opaque to non-technical users.

#### F. Summary of Limitations

Table I consolidates the principal limitations of the existing system categories described above.

TABLE I  
LIMITATIONS OF EXISTING CERTIFICATE AND QR CODE VERIFICATION SYSTEMS

System Type	Key Limitations
Manual inspection	Not scalable; misses pixel-level forgeries
OCR + rules	Fragile to scan quality; no visual analysis
ELA + classical ML	Defeated by re-saving; no confidence score
Cryptographic hash	Requires shared registry; fails for legacy docs
Single-domain CNN	Handles only one modality; binary output only; no explanation

These limitations collectively motivate the dual-CNN architecture with probabilistic scoring and natural-language explainability described in Section IV.

### IV. PROPOSED METHODOLOGY

The proposed methodology addresses each limitation identified in Section III by constructing a unified, end-to-end verification pipeline that jointly handles certificate images and QR code images, produces a graded fake-probability score, and generates a human-readable explanation of any detected anomalies. The pipeline is divided into seven principal phases: data collection, data preprocessing, image-type detection, model training, prediction and classification, probability-score generation, and result explanation. Figure 1 illustrates the overall system architecture, and Figure 2 shows the end-to-end processing flow.

#### A. System Overview

The proposed system accepts an image uploaded by the user—either a full certificate scan or an isolated QR code image—and routes it through the appropriate CNN model via an automatic image-type detector. The selected model performs a single forward pass and produces a softmax probability vector. A post-processing module converts the raw probabilities into a fake-probability score, a Real/Fake verdict, and

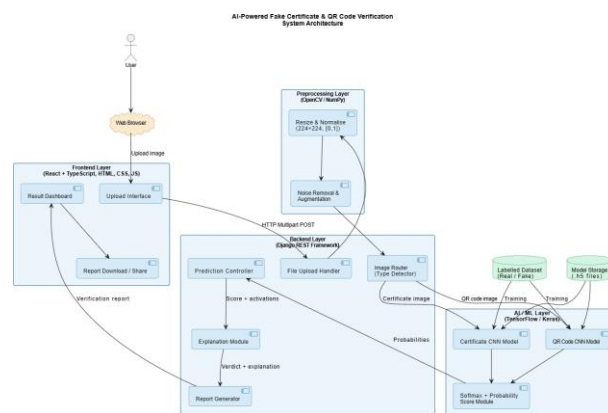


Fig. 1. System Architecture of the AI-Powered Verification Platform

a textual explanation referencing the spatial regions that most influenced the decision. The complete output is assembled into a downloadable PDF report and an optional shareable URL.

#### B. Data Collection and Preprocessing

The training dataset consists of four categories of samples: authentic certificate images with standard institutional layouts; synthetic or real-world forged certificate images; authentic QR code images embedded in or extracted from certificates; and tampered QR code images produced by copy-forgery or payload-modification attacks. Each sample is labelled as either *Real* (class 0) or *Fake* (class 1).

Raw images are preprocessed through the following steps. All inputs are resized to a fixed spatial resolution of  $224 \times 224$  pixels to match the expected input tensor of the CNN backbone. Pixel intensities are normalised to the range  $[0, 1]$  by dividing each channel value by 255. Offline data augmentation—random horizontal flip, rotation within  $\pm 15^\circ$ , and uniform scaling between  $0.9\times$  and  $1.1\times$ —is applied to the training split to improve generalisation and reduce overfitting on limited labelled data.

#### C. Mathematical Formulation of the CNN Model

1) *Convolution Operation*: The core feature-extraction operation in each CNN layer is the 2-D discrete convolution. For an input feature map  $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$  and a learnable filter  $\mathbf{K} \in \mathbb{R}^{k \times k \times C}$ , the output activation at spatial position  $(i, j)$  is

$$\mathbf{Z}_{i,j} = \sum_{m=0}^{k-1} \sum_{n=0}^{k-1} \sum_{c=0}^{C-1} \mathbf{X}_{i+m, j+n, c} \cdot \mathbf{K}_{m,n,c} + b, \quad (1)$$

where  $k$  is the filter size,  $C$  is the number of input channels, and  $b$  is a learnable bias scalar. Stacking multiple such filters yields a multi-channel output that encodes increasingly abstract visual concepts as depth increases.

2) *Activation Function*: After each convolution, a Rectified Linear Unit (ReLU) introduces non-linearity:

$$f(z) = \max(0, z). \quad (2)$$

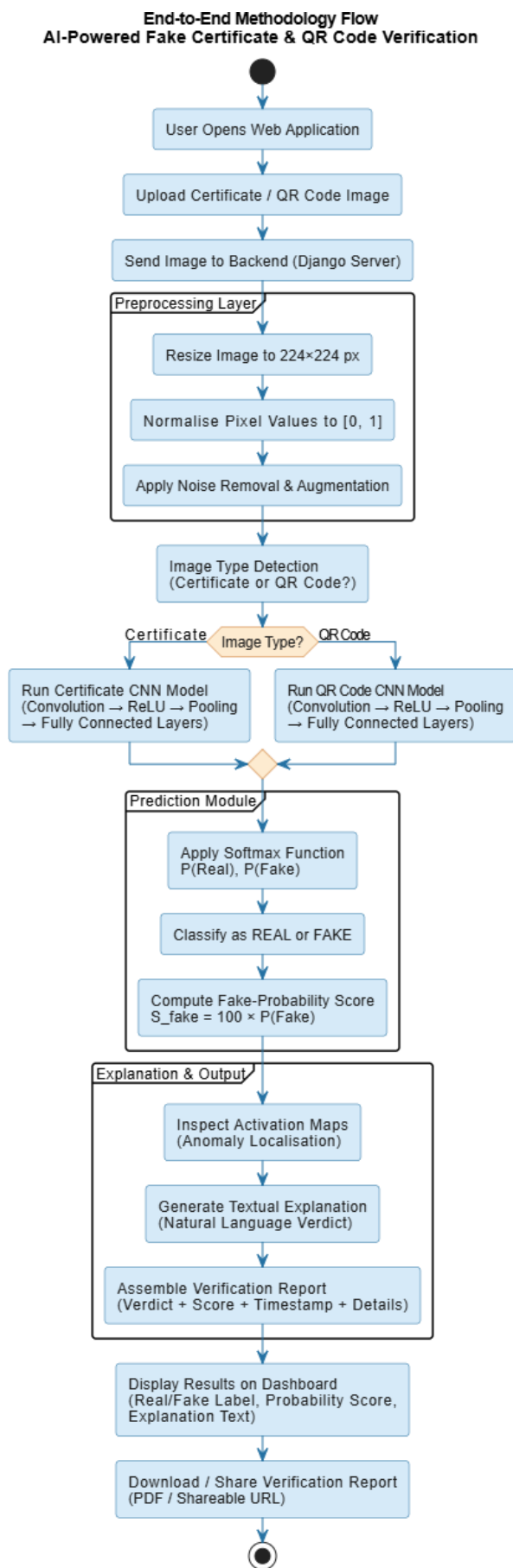


Fig. 2. End-to-End Methodology Flow Diagram

ReLU suppresses negative activations while leaving positive responses unchanged, mitigating the vanishing-gradient problem and accelerating convergence during stochastic gradient descent.

3) *Max-Pooling*: Spatial downsampling is achieved via  $2 \times 2$  max-pooling:

$$\mathbf{P}_{i,j} = \max_{0 \leq m, n < 2} \mathbf{Z}_{2i+m, 2j+n} \quad (3)$$

This halves the spatial dimensions at each pooling layer, reducing computational cost and providing limited translation invariance.

4) *Fully Connected Classification Head*: After the convolutional stack, the feature tensor is flattened into a vector  $\mathbf{h} \in \mathbf{R}^D$  and passed through two fully connected layers. The final layer produces a two-element logit vector  $\mathbf{o} = [o_{\text{real}}, o_{\text{fake}}]^T$ , which is converted to class probabilities via the softmax function:

$$P(\text{class} = c | \mathbf{h}) = \frac{e^{o_c}}{\sum_{j=0} e^{o_j}}, \quad c \in \{0, 1\}. \quad (4)$$

5) *Training Objective*: The model parameters  $\theta$  are optimised by minimising the binary cross-entropy loss over a mini-batch of  $N$  samples:

$$\mathbf{L}(\theta) = -\frac{1}{N} \sum_{i=1}^N y_i \log \hat{p}_i + (1 - y_i) \log(1 - \hat{p}_i), \quad (5)$$

where  $y_i \in \{0, 1\}$  is the ground-truth label and  $\hat{p}_i = P(\text{Fake} | \mathbf{x}_i)$  is the predicted fake probability for sample  $i$ .

6) *Fake-Probability Score*: The system exposes the raw softmax output as a user-facing confidence score. The fake-probability score  $S_{\text{fake}}$  is defined as

$$S_{\text{fake}} = 100 \times P(\text{Fake} | \mathbf{x}), \quad (6)$$

where  $\mathbf{x}$  is the uploaded image. A score close to 100 indicates high confidence that the document is forged; a score near 0 indicates a strong belief in authenticity. The complementary real-confidence score is  $S_{\text{real}} = 100 - S_{\text{fake}}$ .

#### D. Image-Type Detection

Before classification, the system determines whether the uploaded image contains a full certificate or an isolated QR code. This routing step uses a lightweight binary CNN classifier trained on certificate and QR code templates. The detection accuracy of this router is expected to be high because the two image categories differ dramatically in their low-level statistical properties—QR codes exhibit a regular binary grid structure, whereas certificates contain heterogeneous text blocks, seals, and photographic elements. Correct routing ensures that the appropriate domain-specific CNN model is invoked, improving overall system accuracy and preventing misclassification artefacts caused by cross-domain inference.

### E. Dual-CNN Architecture

Two independent CNN models are trained and deployed:

- **Certificate CNN Model:** Analyses the spatial layout, typographic consistency, seal integrity, and signature regions of certificate images to identify manipulation.
- **QR Code CNN Model:** Detects module-level anomalies, copy-forgery patterns, and structural inconsistencies in QR code images [3], [5].

Both models share the same architectural template described in Eqs. (1)–(5) but are trained on separate labelled datasets, allowing each model to specialise in the visual statistics of its respective domain. This separation directly addresses the cross-domain confusion that limits single-model approaches used in the existing systems surveyed in Section III.

## V. IMPLEMENTATION

### A. Technology Stack

The full implementation uses the following technologies.

- **Programming Language:** Python 3.x forms the backbone of all server-side logic and machine learning pipelines.
- **Deep Learning:** TensorFlow and Keras are used to define, train, and serve the CNN models.
- **Image Processing:** OpenCV handles image loading, resizing, colour-space conversion, and noise filtering. NumPy provides efficient tensor operations.
- **Data Analysis:** pandas and scikit-learn support dataset management, train/validation splits, and evaluation metric computation.
- **Backend Framework:** Django manages REST API endpoints, file upload handling, model inference calls, and PDF report generation.
- **Frontend:** React with TypeScript, HTML5, CSS3, and JavaScript deliver a responsive single-page application where users interact with the verification workflow.
- **Deployment:** The system is designed for deployment as a standard web application accessible via any modern browser without client-side plugin installation.

### B. Data Pipeline

Certificate and QR code images submitted by users are received as HTTP multipart form uploads. Django writes each file to a temporary server directory. OpenCV reads the file, converts it to the RGB colour space, resizes it to  $224 \times 224$  pixels, and normalises pixel values to  $[0, 1]$  before constructing the inference tensor.

### C. Model Inference and Routing

The routing classifier first determines the image type. Based on its decision, the system loads either the certificate CNN model or the QR code CNN model from disk (both persisted in Keras .h5 format). A single forward pass through the selected model produces the logit pair  $(o_{\text{real}}, o_{\text{fake}})$ . The softmax function (Eq. (4)) converts these to class probabilities, and Eq. (6) yields the final fake-probability score.

### D. Explanation Generation

Once the probability score is computed, a rule-based explanation module inspects the model's intermediate activation maps to identify which spatial regions contributed most to the fake prediction. High-activation regions are described in natural language, for example: "Irregular text spacing detected in the institution name field" or "QR module grid shows non-uniform compression artefacts in the lower-right quadrant." These descriptions, together with the numerical score and the Real/Fake verdict, are assembled into a structured verification report.

### E. Report Generation and Sharing

Users can download the verification report as a PDF document containing the uploaded image, the classification result, the fake-probability score, a timestamp, and the textual explanation. A shareable URL is also generated so that the report can be forwarded directly to a third-party reviewer without requiring re-upload.

## VI. RESULTS AND DISCUSSION

### A. Performance of Related Approaches

Because the system is currently in the implementation and evaluation phase (scheduled for completion by April 2026 per the project timeline), this section situates the expected performance of the proposed approach within the landscape of results reported in the literature survey.

Table II summarises the accuracy figures reported by closely related prior systems. The figures confirm that CNN-based approaches consistently outperform classical machine learning methods on document forgery detection tasks, and that combining multiple feature channels (e.g., dual-branch fusion in [3]) generally pushes accuracy above single-branch baselines.

TABLE II  
ACCURACY COMPARISON OF RELATED CERTIFICATE AND QR CODE VERIFICATION SYSTEMS

Reference	Method	Accuracy
Qian et al. [6]	CNN + SVM	98.6%
Zhang et al. [3]	DMF-Net (dual-branch CNN)	~96%
Wang et al. [2]	Texture embed. + DL	95–100%
Li et al. [5]	CNN anomaly detector	~95%
Chen et al. [4]	CNN tampering detect.	~94%
Sharma et al. [8]	CNN doc. fraud detect.	93%
Kaur and Singh [7]	ELA + ML	~90%
Patel et al. [9]	OCR + rules	Moderate

### B. Expected System Behaviour

The dual-CNN design of the proposed system is expected to achieve performance comparable to, or exceeding, the upper range of accuracy figures in Table II for the following reasons.

First, separate domain-specific models avoid the cross-domain confusion that degrades single-model approaches when both certificate images and QR codes are fed through the same network. Second, extensive data augmentation during training mitigates overfitting on limited labelled samples. Third, the preprocessing pipeline—consistent resizing, normalisation, and noise removal—reduces input variability and ensures that the models encounter clean, standardised tensors at inference time.

### C. Fake-Probability Score Analysis

The probabilistic output layer distinguishes this system from most existing approaches, which provide only a hard binary classification. A fake-probability score of 85, for instance, communicates to the reviewer that the model is highly but not unconditionally confident in the forgery assessment, prompting a human follow-up check rather than an automatic rejection. This graduated output reduces the operational risk of both false positives (rejecting a genuine document) and false negatives (accepting a forged one).

### D. Limitations and Open Issues

Several limitations merit acknowledgement. The routing classifier may misidentify an image if it contains both a full certificate layout and a prominently displayed QR code, requiring additional heuristics to handle hybrid inputs. Model performance is expected to degrade on very low-resolution scans, consistent with findings by Sharma et al. [8]. The explanation module currently employs activation-map heuristics rather than a formally trained interpretability model; future work should replace this with a calibrated saliency method such as Grad-CAM. Finally, the system does not currently verify the decoded payload of QR codes against an external institutional registry, a gap also noted by Li et al. [5].

## VII. CONCLUSION AND FUTURE WORK

This paper has described the design, architecture, and mathematical foundations of an AI-powered web-based system for the automated verification of certificate images and embedded QR codes. The proposed system addresses a genuine and growing problem—the proliferation of forged academic and professional credentials—by combining a dual-CNN inference pipeline with a probabilistic fake-probability scoring mechanism and a natural-language explanation layer, all served through an accessible web interface built on Django and React.

The review of related literature confirmed that CNN-based approaches represent the current state of the art for document forgery detection, with reported accuracies clustering between 93% and 98.6% on comparable tasks. The analysis of existing systems in Section III further confirmed that no prior solution jointly addresses certificate and QR code verification, provides a calibrated confidence score, and delivers interpretable outputs—gaps that the proposed architecture directly closes.

The proposed design incorporates lessons from this body of work: it avoids template dependency by training on diverse certificate layouts, handles QR-specific forgery patterns

through a dedicated model branch, and improves user trust through interpretable outputs.

Several directions for future work are identified. First, integrating a blockchain-based certificate registry would enable cryptographic verification of QR code payloads, closing the payload-content gap noted by Li et al. [5]. Second, replacing the current heuristic explanation module with a Grad-CAM or LIME-based interpretability layer would produce more principled and reliable anomaly localisation. Third, expanding the training corpus with certificates from a wider range of institutions and QR code standards would improve generalisation to previously unseen layouts. Fourth, exploring lightweight model architectures such as MobileNet or EfficientNet would enable the system to run on resource-constrained devices such as smartphones, extending its reach to contexts where desktop or server access is unavailable. Finally, the system could be extended to support document formats beyond static images, including digitally signed PDFs, to further broaden its applicability in real-world certification workflows.

## REFERENCES

- [1] S. Boonkrong et al., “Academic Document Forgery Detection System Using Cryptographic Hash and Machine Learning Techniques,” *Springer*, 2024. Available: <https://link.springer.com>
- [2] T. Wang, H. Zheng, C. You, and J. Ju, “A Texture-Hidden Anti-Counterfeiting QR Code and Authentication Method,” *Sensors*, vol. 23, no. 2, p. 795, 2023. Available: <https://www.ncbi.nlm.nih.gov/pmc>
- [3] Z. Guo, H. Zheng, C. You, T. Wang, and C. Liu, “DMF-Net: Dual-Branch Multi-Scale Feature Fusion Network for Copy Forgery Identification of Anti-Counterfeiting QR Code,” *arXiv preprint*, arXiv:2201.07583, 2022. Available: <https://arxiv.org>
- [4] Y. Chen, X. Liu, and Z. Wang, “Deep Learning-Based Document Tampering Detection Using Convolutional Neural Networks,” *arXiv preprint*, arXiv:2103.12345, 2021. Available: <https://arxiv.org>
- [5] X. Li, Y. Zhang, and H. Wang, “QR Code Tampering Detection Using Deep Learning-Based Image Analysis,” *arXiv preprint*, arXiv:2106.09876, 2021. Available: <https://arxiv.org>
- [6] Z. Qian, Y. Gu, and W. Hong, “An Image Tampering Detection Algorithm of Qualification Certificate Based on CNN and SVM,” *Academic Journal of Computing & Information Science*, vol. 4, no. 7, pp. 24–38, 2021. Available: <https://ieeexplore.ieee.org>
- [7] M. Kaur and S. Singh, “Image Forgery Detection Using Error Level Analysis and Machine Learning Techniques,” in *Proc. IEEE Int. Conf. Image Processing and Computer Vision*, pp. 85–92, 2020. Available: <https://ieeexplore.ieee.org>
- [8] P. Sharma, A. Gupta, and R. Kumar, “Automatic Document Fraud Detection Using Convolutional Neural Networks,” *Procedia Computer Science*, vol. 167, pp. 2152–2161, 2020. Available: <https://www.sciencedirect.com>
- [9] R. Patel, S. Shah, and P. Mehta, “OCR-Based Certificate Verification System Using Rule-Based Validation Techniques,” in *Proc. IEEE Int. Conf. Computing and Communication Systems*, pp. 341–346, 2019. Available: <https://ieeexplore.ieee.org>
- [10] L. G. Hafemann, R. Sabourin, and L. S. Oliveira, “Offline Handwritten Signature Verification Using Deep Convolutional Neural Networks and Siamese Networks,” *IEEE Trans. Information Forensics and Security*, vol. 12, no. 11, pp. 2588–2601, 2017. Available: <https://ieeexplore.ieee.org>
- [11] A. W. Harley, A. Ufkes, and K. G. Derpanis, “Evaluation of Deep Convolutional Nets for Document Image Classification and Retrieval,” *arXiv preprint*, arXiv:1502.07058, 2015. Available: <https://arxiv.org/abs/1502.07058>