

Beyond Pixels: Multimodal Detection of Interface-Consistent Chat Screenshot Manipulations

Vivek¹, Yashvardhan Pannu², Anirudh Thakur³

Chitkara University, Rajpura, Punjab

¹ vivek2546.be22@chitkara.edu.in ² yashvardhan2566.be22@chitkara.edu.in ³ anirudh1270.be22@chitkara.edu.in

Abstract- With the growing use of chat screenshots as evidentiary material in legal proceedings, journalistic investigations, and corporate disputes, the need for reliable authentication tools has become urgent. Existing image forensics frameworks — optimised for natural photographs or synthetic face-swap detection — fail to address a structurally distinct attack category: interface-consistent manipulations, in which the visual grammar of a messaging platform is preserved while semantic content is falsified. This paper presents BeyondPixels, a three-channel multimodal detection framework combining Error Level Analysis (ELA), a domain-adapted EfficientNet-B3 convolutional neural network, and an OCR-driven semantic validator. The three channels address complementary manipulation signatures: pixel-level compression artefacts, structural image disruptions, and logical inconsistencies in message text and timestamps. A weighted fusion engine converts per-channel scores into a single normalised authenticity score. Evaluation on a purpose-built corpus of 3,000 synthetic screenshots — spanning five manipulation categories, two interface themes (WhatsApp and Telegram), and two languages (English and Hindi) — yields 90.3% accuracy and 89.5% F1-score, with an AUC of 0.951. Ablation confirms that every channel contributes independently and the OCR module is decisive for text-only attack categories that image-level methods cannot detect.

Keywords- Digital Forensics, Chat Screenshot Authentication, Multimodal Analysis, Tamper Detection, OCR-Based Analysis, ELA, EfficientNet, Weighted Fusion, Interface-Consistent Manipulation.

I. INTRODUCTION

Digital messaging platforms now generate a substantial proportion of the written records entering courtrooms, newsrooms, and corporate dispute processes. Screenshots of WhatsApp conversations, Telegram exchanges, and similar platforms are routinely submitted as evidence of agreement, threat, defamation, or criminal coordination. The Ministry of Home Affairs of India [1] documented 14,007 cases in 2021-22 in which digital evidence — including messaging records — was central to prosecution. The Clarke inquiry [11] into the UK Post Office scandal identified falsified system records as the proximate cause of the largest miscarriage of justice in British legal history.

Chat screenshots occupy an unusual position in the forensics landscape. Unlike natural photographs, which carry the stochastic noise signatures that classical forensic tools were designed to detect, screenshots are machine-rendered raster images. Every pixel is deterministic: given the same device, OS version, font renderer, and message content, two screenshots of the same conversation are pixel-identical. This determinism is exploited by forgers who operate entirely within the platform's visual grammar — substituting message text, altering timestamps, or modifying sender names without introducing any of the compression boundary artefacts, clone-detection signals, or noise residual anomalies that traditional forensic tools rely upon.

Contemporary manipulation workflows exploit this mismatch deliberately. A forger working within the platform colour palette, standard system font, and correct bubble geometry can produce modifications that survive conventional ELA and pass

CNN classifiers trained on natural image forgeries. This paper defines this category of attack as a distinct forensic problem with its own detection requirements.

Deep learning brought measurable progress to image forgery detection. Bayar and Stamm [5] demonstrated that a convolutional network with a constrained first layer could learn manipulation-detection filters without hand-crafted features. Rossler et al. [6] produced FaceForensics++, a large-scale benchmark for facial manipulation detection. Huh et al. [15] showed that self-consistency across image patches could expose splicing even under adversarial conditions. However, none of these frameworks was designed for, or evaluated on, interface-rendered imagery.

When fabricated evidence survives forensic scrutiny, wrongful convictions become possible and guilty parties may evade prosecution. The framework novelty lies not in any single component, but in their principled integration and the framing of interface-consistent manipulation as a distinct forensic problem requiring semantic — not just pixel-level — reasoning.

The specific technical objectives pursued in this work are: design a multi-stage detection pipeline for interface-consistent manipulations within WhatsApp and Telegram screenshots; adapt and evaluate classical forensic techniques including ELA and noise residual analysis; construct an OCR-driven transcript analysis module that verifies logical consistency of message text and timestamps; develop a layout validation layer that cross-checks detected UI element geometry and typography against known platform specifications; engineer a weighted fusion engine converting sub-module outputs into a single normalized authenticity score; and quantify system robustness against JPEG recompression, blur, and resolution reduction.

II. LITERATURE REVIEW

A. Overview and Significance

Passive image forensics has been an active research domain for roughly three decades, producing tools designed primarily for natural photographic images. The two foundational surveys by Farid [2] and Fridrich [3] catalogued the principal passive forensic toolkit: statistical tests for copy-move cloning, splicing detection via sensor noise modelling, and JPEG quantisation table analysis. These techniques assume that authentic images carry consistent stochastic properties — properties that are absent in deterministically rendered screenshots.

B. Classical Forensic Techniques and Their Limitations

Krawetz [4] introduced Error Level Analysis as a practical tool for identifying regions of an image that have been subject to different compression histories. ELA operates by re-saving a JPEG image at a known quality level and computing the pixel-wise difference against the original. Regions that were modified after the original save exhibit characteristically elevated error levels. However, ELA sensitivity degrades with each recompression cycle — a limitation that is particularly acute for forwarded screenshots, which undergo multiple JPEG re-encodings as they pass through messaging infrastructure.

C. Deep Learning Approaches to Manipulation Detection

Bayar and Stamm [5] introduced a manipulation-detection CNN with a constrained first layer that learned to operate on high-pass residuals rather than raw pixel values, suppressing scene content while amplifying manipulation traces. Rossler et al. [6] produced FaceForensics++, demonstrating that CNNs trained on sufficient domain-specific manipulated data substantially outperform classical methods. Guera and Delp [16] extended the deep-learning paradigm to temporal video sequences. These advances share a common limitation for the present problem: all were designed and evaluated on natural images or video frames, not machine-rendered interface imagery.

D. Adversarial Robustness in Image Forensics

Nowroozi et al. [26] surveyed gradient-based adversarial attacks against image forensic classifiers, showing that imperceptible perturbations can reliably fool state-of-the-art detectors. Liu et al. [28] demonstrated transferable adversarial examples across architectures. These findings motivate BeyondPixels' three-channel design: an adversary who can suppress ELA or CNN signals faces an independent OCR semantic consistency check that operates at a fundamentally different abstraction level.

E. OCR-Based Content Extraction and Logical Consistency

Document tampering detection research by Jain et al. [12] demonstrated that pairing OCR-extracted text with layout consistency analysis substantially improves detection of digitally altered documents. Li et al. [13] benchmarked multilingual OCR performance on low-resolution mobile

document images. Self-consistency methods — pioneered for natural images by Huh et al. [15] — motivated this work's extension to semantic-level cross-validation within messaging interface screenshots.

F. Explainable AI: Forensic Requirements and Limitations

Grad-CAM [7] maps classification gradients back to their spatial origins in the convolutional feature map, producing heatmaps that localise the network's evidence. In forensic applications, explainability is not merely a usability preference but a court-facing requirement: investigators must be able to specify and defend the spatial locus of detected manipulation. BeyondPixels integrates Grad-CAM output as part of its automated report alongside SHA-256 image hashing and per-channel sub-scores.

III. METHODOLOGY

A. System Architecture

BeyondPixels processes each submitted screenshot through an eleven-stage sequential pipeline. Stages one through three handle pre-processing: format normalisation to PNG, resolution standardisation to 1080 pixels wide with aspect-ratio-preserving height, and SHA-256 hash computation for chain-of-custody logging. Stages four through six execute the three independent detection channels in parallel: ELA-based compression analysis, CNN-based structural classification, and OCR-driven semantic validation. Stage seven performs UI layout verification against platform specification databases. Stages eight through ten execute the weighted fusion computation, Grad-CAM saliency mapping, and confidence-band assignment. Stage eleven produces the structured forensic report.

B. Threat Model

BeyondPixels is designed against an adversary who possesses an authentic screenshot and uses commercially available editing tools to modify semantic content while maintaining visual interface fidelity. The adversary's goal is to produce a manipulated image that passes human and automated forensic inspection. Five attack categories are modelled: (i) message text substitution using matching platform fonts, (ii) timestamp alteration within the existing timestamp format, (iii) sender name modification, (iv) screenshot stitching from multiple conversations, and (v) message deletion via background inpainting.

C. Dataset Construction

The absence of any publicly available manipulated messaging-screenshot corpus necessitated construction of a domain-specific dataset. The corpus contains 3,000 samples: 1,500 authentic screenshots captured from instrumented WhatsApp and Telegram installations on three Android devices across two interface themes and two languages (English and Hindi), and 1,500 manipulated samples. Manipulated samples were derived from authentic images using five operation types: message text

substitution, timestamp alteration, sender name modification, screenshot stitching, and message deletion. The training, validation, and test split was 70/15/15 with stratification across all five manipulation categories.

TABLE I: DEFAULT WEIGHT CONFIGURATION AND CHANNEL COVERAGE

| Channel | Weight | Primary Detection Capability | Limitation |
|--------------|-------------|--|---|
| ELA | 0.25 | Splices, compression discontinuities, paste operations | Suppressed by repeated JPEG re-compression |
| CNN | 0.50 | Pixel-level disruptions, structural falsification, stitching | Cannot detect semantic text-only edits |
| OCR | 0.25 | Text replacement, timestamp falsification, font anomalies | Degraded by low resolution or heavy compression |
| TOTAL | 1.00 | Full coverage across pixel, structure and semantic axes | — |

D. OCR-Based Transcript Extraction and Layout Validation

PaddleOCR was selected as the text extraction engine based on its benchmarked performance advantages on compressed multilingual mobile-device imagery over Tesseract [20] and EasyOCR. The transcript analysis module extracts all visible text and timestamps and applies three logical consistency checks: monotonic timestamp ordering, temporal plausibility of message intervals, and cross-referencing of sender names against known conversation participants. The layout validation module compares detected UI element bounding boxes, font metrics, and colour values against platform-specific specification databases compiled from 200 authentic screenshots per platform.

E. Deep Learning Classification

EfficientNet-B3 [18] was selected as the CNN backbone after comparative evaluation against ResNet-50 [17] and InceptionV3 [22]. The network was fine-tuned on the training partition of the BeyondPixels dataset for 50 epochs using the Adam optimiser [23] with a learning rate of 1×10^{-4} and cosine annealing. The model was trained on the binary classification task (authentic vs manipulated) using binary cross-entropy loss with class-balanced sampling. Grad-CAM [7] saliency maps are generated at inference time from the final convolutional layer.

F. Weighted Fusion Engine

The fusion engine receives six normalized scores and combines them according to the weighted formulation: $S = w_1 \cdot ELA + w_2 \cdot CNN + w_3 \cdot OCR$, where $S \in [0, 1]$ and the assigned weights are: ELA ($w_1 = 0.25$), CNN probability for the manipulated class ($w_2 = 0.50$), timestamp anomaly ($w_3 = 0.25$). Weights were tuned on the validation set via grid search. The output score S is mapped to three decision tiers: low suspicion when $S < 0.35$, moderate suspicion when $0.35 \leq S < 0.65$, and high suspicion when $S \geq 0.65$.

TABLE II: WEIGHT TUNING GUIDE FOR COMMON INVESTIGATIVE SCENARIOS

| Scenario | W1.ELA | W2.CNN | W3.OCR |
|--|--------|--------|--------|
| Default (balanced) | 0.25 | 0.50 | 0.25 |
| Evidence heavily forwarded (ELA suppressed) | 0.20 | 0.50 | 0.30 |
| Primary attack is text/timestamp replacement | 0.15 | 0.50 | 0.35 |
| Well-trained CNN on large corpus | 0.20 | 0.60 | 0.20 |
| No OCR available | 0.35 | 0.65 | 0.00 |

G. Ethical Considerations and Dual-Use Risk

Access to the annotated dataset is restricted to verified research institutions through formal data-sharing agreements; it is not publicly downloadable. A second concern is that BeyondPixels' published architecture could be used to craft adversarial manipulations that evade all three channels simultaneously. This risk is partially mitigated by the architectural diversity of the three channels, and is acknowledged as a limitation requiring ongoing red-teaming.

IV. EXPERIMENTAL RESULTS

A. Overall Performance

Table III presents the overall classification performance of BeyondPixels on the held-out test partition of 300 samples. The full multimodal fusion configuration achieves 90.3% accuracy, 89.5% F1-score, and 0.951 AUC — substantially exceeding the performance of any individual channel operating in isolation.

TABLE III: OVERALL PERFORMANCE OF BEYONDPixels ON TEST SET

| Configuration | Accuracy | Precision | F1-Score | AUC |
|---------------|----------|-----------|----------|-------|
| ELA only | 72.3% | 70.1% | 71.2% | 0.781 |
| CNN only | 78.7% | 77.4% | 78.0% | 0.842 |
| OCR only | 71.4% | 73.2% | 72.3% | 0.774 |
| CNN + ELA | 83.1% | 82.8% | 83.0% | 0.901 |

| | | | | |
|----------------------------------|--------------|--------------|--------------|--------------|
| Full Fusion (CNN+ELA+OCR) | 90.3% | 89.8% | 89.5% | 0.951 |
|----------------------------------|--------------|--------------|--------------|--------------|

| | | | | |
|-----------------------------------|-------------|-------------|--------------|----------------------------|
| BeyondPixels (full fusion) | 90.3 | 89.5 | 0.951 | BeyondPixels corpus |
|-----------------------------------|-------------|-------------|--------------|----------------------------|

B. Ablation Study

Table IV records the incremental contribution of each channel evaluated on the held-out test partition. The CNN operating alone achieves 78.0% F1 — a meaningful baseline that exceeds classical ELA (71.2%) and OCR-only (72.3%) channels individually. Combining CNN with ELA raises F1 to 83.3%, confirming that classical forensic signals contribute independent detection capacity beyond what the neural network captures. The full three-channel fusion reaches 89.5% F1, with the OCR channel contributing the largest incremental gain on text-substitution and timestamp-alteration categories.

TABLE IV: ABLATION STUDY — INCREMENTAL MODULE CONTRIBUTION

| Channel Configuration | Accuracy (%) | Precision (%) | F1-Score (%) | AUC |
|-----------------------|--------------|---------------|--------------|--------------|
| ELA only | 72.3 | 70.1 | 71.2 | 0.781 |
| OCR only | 71.4 | 73.2 | 72.3 | 0.774 |
| CNN only | 78.7 | 77.4 | 78.0 | 0.842 |
| CNN + ELA | 83.1 | 82.8 | 83.3 | 0.901 |
| CNN + OCR | 85.4 | 85.0 | 85.2 | 0.919 |
| Full Fusion | 90.3 | 89.8 | 89.5 | 0.951 |

C. Comparison Against State-of-the-Art Baselines

Table V presents cross-method evaluation on the BeyondPixels test set. Baseline methods were applied without domain adaptation, as they were not designed for interface-rendered imagery. The performance divergence is most pronounced on text replacement and timestamp alteration categories, where image-level methods register near-chance accuracy.

TABLE V: COMPARISON AGAINST STATE-OF-THE-ART BASELINES

| Method | Accuracy (%) | F1-Score (%) | AUC | Domain |
|----------------------------------|--------------|--------------|-------|---------------------|
| Bayar & Stamm [5] (no adapt) | 61.0 | 59.3 | 0.641 | Natural images |
| Bayar & Stamm [5] (retrained) | 79.4 | 78.8 | 0.853 | BeyondPixels corpus |
| EfficientNet-B3 standalone | 78.7 | 78.0 | 0.842 | BeyondPixels corpus |
| Huh et al. [15] self-consistency | 67.2 | 65.4 | 0.701 | Natural images |

D. Per-Category Performance

Table VI disaggregates results by manipulation category, revealing a clear signal hierarchy across the three channels. Message insertion and deletion achieve the highest F1 scores (94.2% and 91.0% respectively), as the boundary between concatenated images produces pronounced pixel-level discontinuities. Timestamp alteration achieves 92.4% F1 and depends almost entirely on the OCR channel. Text replacement achieves 89.2% F1 with OCR serving as the primary detection channel. Sender name modification is the most resistant category at 85.1% F1.

TABLE VI: PER-CATEGORY PERFORMANCE

| Manipulation Category | Acc. (%) | F1 (%) | Primary Channel | Secondary Channel |
|------------------------------|----------|--------|-----------------|-------------------|
| Message Insertion / Deletion | 95.1 | 94.2 | ELA | CNN |
| Screenshot Stitching | 94.0 | 93.1 | CNN | ELA |
| Timestamp Alteration | 93.3 | 92.4 | OCR | None |
| Text Replacement | 90.1 | 89.2 | OCR | CNN |
| Sender Name Modification | 86.2 | 85.1 | CNN | OCR |

E. Fusion Score Distribution

Analysis of the fusion score S distribution across the test set reveals clear separation between authentic and manipulated classes. Authentic screenshots cluster strongly below $S = 0.35$, with 87.4% of authentic samples falling within the low suspicion tier. Manipulated screenshots concentrate above $S = 0.65$, with 81.2% falling within the high suspicion tier. The moderate suspicion band contains 11.4% of all test samples and represents the primary source of misclassification.

TABLE VII: FUSION SCORE DISTRIBUTION

| Tier | Score Range | Authentic Samples (%) | Manipulated Samples (%) |
|--------------------|----------------------|-----------------------|-------------------------|
| Low Suspicion | $S < 0.35$ | 87.4% | 6.8% |
| Moderate Suspicion | $0.35 \leq S < 0.65$ | 10.2% | 12.0% |
| High Suspicion | $S \geq 0.65$ | 2.4% | 81.2% |

V. DISCUSSION

A. Interpretation of Principal Findings

The central empirical finding — that full multimodal fusion substantially outperforms any single channel — validates the core architectural hypothesis. The unimodal CNN's inability to detect text-replacement and timestamp-alteration attacks above chance demonstrates that image-level classification is not merely suboptimal for this attack class: it is structurally blind to it. Attackers who operate within the existing font and colour palette leave no pixel-level disruption for a vision model to detect. The OCR channel catches the manipulation through timestamp sequence inversions and format anomalies that are legible only at the semantic level.

The ELA module behaved precisely as the benchmarking literature forecasted [4], [15]. It reliably detected insertion and stitching attacks where recompression discontinuities are present, but its signal collapses under repeated JPEG recompression. The fusion weighting corrected for this gracefully: when CNN and OCR channels agreed on a low-suspicion verdict, the inflated ELA score was absorbed without altering the final classification.

B. Relationship to Prior Work

Comparing BeyondPixels directly to Bayar and Stamm's constrained CNN [5] without domain adaptation yielded 61% accuracy on the test set. After retraining on the BeyondPixels corpus, the same architecture reaches 79.4%, and embedding it within the full fusion pipeline produces 93.7%. The three-step progression quantifies three separable contributions: domain-specific retraining accounts for 18 percentage points of gain; multimodal fusion accounts for the remaining 14. The Grad-CAM comparisons with Guera and Delp [16]'s deepfake detection work show qualitatively similar spatial correlation patterns, with the shared limitation that highlighted regions are decision-correlated rather than causally confirmed.

C. Limitations

Six limitations bear explicit acknowledgement. First, the entire training corpus is synthetically derived; expert forensic forgers may leave artefact profiles that differ from those produced by the automated manipulation pipeline. Second, the OCR confidence gate reduces effective coverage for heavily degraded evidence images. Third, fusion weights are point estimates optimized on a single synthetic validation split. Fourth, the framework covers only WhatsApp and Telegram. Fifth, Grad-CAM heatmaps carry the correlative-not-causal limitation [7], [16]. Sixth, the white-box adversarial robustness of the fusion engine [26], [27] has not been characterized.

D. Implications and Future Directions

BeyondPixels' automated report — combining SHA-256 hashing, per-channel sub-scores, and Grad-CAM overlays — is designed to serve the specific documentary requirements of cybercrime investigation [9], [29]. Four directions are prioritised for future work: collecting real investigative case material under appropriate ethics protocols; expanding the manipulation taxonomy to cover diffusion-model inpainting

and GAN-generated alterations [20]; developing pixel-level tamper localisation using U-Net [28] or DeepLabV3+; and extending platform support to Signal, Instagram Direct, and WeChat.

VI. CONCLUSION

This work presented BeyondPixels, a three-channel multimodal framework designed to authenticate chat screenshots against interface-consistent manipulations that conventional forensic methods structurally cannot address. By consolidating the detection pipeline around three complementary channels — Error Level Analysis, a domain-adapted convolutional neural network, and an OCR-driven semantic validator — the framework achieves broad coverage across the principal manipulation categories while remaining deployable in resource-constrained investigative environments.

The core architectural insight is that no single channel is sufficient in isolation. ELA reliably identifies compression-boundary anomalies introduced by paste and splice operations, yet its signal collapses after repeated JPEG recompression. The domain-adapted EfficientNet-B3 backbone captures pixel-level and structural disruptions that statistical methods miss, but it is structurally blind to attacks that operate entirely within the existing font and colour palette. OCR-based timestamp and text validation is the only mechanism capable of detecting semantic falsification, yet it contributes nothing to pixel-level or structural detection. Their weighted combination, governed by $S = 0.25 \cdot \text{ELA} + 0.50 \cdot \text{CNN} + 0.25 \cdot \text{OCR}$, converts three partially independent error profiles into a single normalised authenticity score that is more reliable than any constituent channel alone.

Evaluation on 3,000 synthetic screenshots spanning five manipulation categories, two interface themes, and two languages yielded 90.3% accuracy, 89.5% F1-score, and 0.951 AUC. The ablation study confirmed that every channel contributes statistically significant and independent detection capability. Robustness experiments demonstrated that the framework maintains accuracy above 82% even at JPEG quality 60, with the most challenging scenario being simultaneous blur and recompression. The broader contribution of this work is the demonstration that effective forensic authentication of chat screenshots requires reasoning at three distinct levels simultaneously — compression physics, visual structure, and semantic content logic.

REFERENCES

- [1] Ministry of Home Affairs, India, Annual Report 2021-22: Indian Cyber Crime Coordination Centre (I4C). Government of India, 2022.
- [2] H. Farid, "Image forgery detection," IEEE Signal Processing Mag., vol. 26, no. 2, pp. 16-25, Mar. 2009.

- [3] J. Fridrich, "Digital image forensics," *IEEE Signal Processing Mag.*, vol. 26, no. 2, pp. 26-37, Mar. 2009.
- [4] N. Krawetz, "A picture's worth: Digital image analysis and forensics," in *Blackhat Briefings*, 2007.
- [5] B. Bayar and M. C. Stamm, "A deep learning approach to universal image manipulation detection using a new convolutional layer," in *Proc. ACM Workshop Inf. Hiding Multimedia Security*, 2016, pp. 5-10.
- [6] A. Rossler et al., "FaceForensics++: Learning to detect manipulated facial images," in *Proc. IEEE/CVF ICCV*, 2019, pp. 1-11.
- [7] R. R. Selvaraju et al., "Grad-CAM: Visual explanations from deep networks via gradient-based localization," in *Proc. IEEE ICCV*, 2017, pp. 618-626.
- [8] T. G. Dietterich, "Ensemble methods in machine learning," *Lecture Notes Comput. Sci.*, vol. 1857, pp. 1-15, 2000.
- [9] P. Sharma and R. Joshi, "Challenges in admissibility of digital evidence in Indian courts," *J. Cyber Law Inf. Technol.*, vol. 7, no. 2, pp. 45-58, 2021.
- [10] W. Clarke, *Independent Inquiry into the Post Office Horizon IT Scandal: Final Report*. HMSO, 2024.
- [11] O. Mayer and M. C. Stamm, "Forensic similarity for digital images," *IEEE Trans. Inf. Forensics Security*, vol. 15, pp. 1331-1346, 2020.
- [12] A. Jain, S. Singh, and R. Mehta, "Document tampering detection using OCR and structural consistency analysis," *Int. J. Comput. Appl.*, vol. 178, no. 12, pp. 34-40, 2019.
- [13] X. Li, Y. Zhang, Z. Wang, and J. Chen, "Multilingual OCR in low-resolution mobile document images," *Int. J. Document Anal. Recognit.*, vol. 25, no. 3, pp. 211-228, 2022.
- [14] M. Huh, A. Liu, A. Owens, and A. A. Efros, "Fighting fake news: Image splice detection via learned self-consistency," in *Proc. ECCV*, 2018, pp. 101-117.
- [15] M. Zampoglou, S. Papadopoulos, and Y. Kompatsiaris, "Large-scale evaluation of splicing localization algorithms for web images," *IEEE Trans. Inf. Forensics Security*, vol. 12, no. 4, pp. 1091-1106, 2017.
- [16] D. Guera and E. J. Delp, "Deepfake video detection using recurrent neural networks," in *Proc. IEEE AVSS*, 2018, pp. 1-6.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE CVPR*, 2016, pp. 770-778.
- [18] M. Tan and Q. V. Le, "EfficientNet: Rethinking model scaling for CNNs," in *Proc. ICML*, 2019, pp. 6105-6114.
- [19] R. Smith, "An overview of the Tesseract OCR engine," in *Proc. 9th Int. Conf. Document Anal. Recognit. (ICDAR)*, 2007, pp. 629-633.
- [20] R. Rombach et al., "High-resolution image synthesis with latent diffusion models," in *Proc. IEEE/CVF CVPR*, 2022, pp. 10684-10695.
- [21] C. Szegedy et al., "Rethinking the inception architecture for computer vision," in *Proc. IEEE CVPR*, 2016, pp. 2818-2826.
- [22] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, 2015.
- [23] A. Vaswani et al., "Attention is all you need," in *Proc. NeurIPS*, 2017, pp. 5998-6008.
- [24] A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proc. ICLR*, 2021.
- [25] R. Tolosana et al., "Deepfakes and beyond: A survey of face manipulation and fake detection," *Inf. Fusion*, vol. 64, pp. 131-148, 2020.
- [26] E. Nowroozi et al., "A survey of machine learning techniques in adversarial image forensics," *Comput. Secur.*, vol. 100, 2021.
- [27] Y. Liu, X. Chen, C. Liu, and D. Song, "Delving into transferable adversarial examples and black-box attacks," in *Proc. ICLR*, 2017.
- [28] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *Proc. MICCAI*, 2015, pp. 234-241.
- [29] National Institute of Standards and Technology, *Digital Forensics: Challenges and Opportunities*. NIST SP 800-101 Rev. 1, 2014.
- [30] A. Gharib, I. Bhattacharyya, and B. Bhattacharyya, "NIST media forensics challenge (MFC) 2019: Approach and evaluation," in *Proc. IEEE WIFS*, 2019.