# A NOVEL DUAL-FRAMEWORK FOR AI SYNTHETIC MEDIA DETECTION BASED ON PHYSIOLOGICAL AND LINGUISTIC INCONSISTENCIES

## Festo K. Magembe[1], Mrindoko Nicholaus [2]

[1] Department of Computer Science and Engineering Mbeya University of Science and Technology, Mbeya, Tanzania

[2] Department of Computer Science and Engineering Mbeya University of Science and Technology, Mbeya, Tanzania

Corresponding Authors E-mail: nicholausmrindoko@gmail.com

Corresponding Authors Mobile No: +255 714 427 241

---

## Abstract

The rapid advancement of generative artificial intelligence has intensified the challenge of detecting AI-synthesized facial media, commonly known as deepfakes. This study introduces a novel dual-framework that fuses physiological and linguistic inconsistency analysis for robust synthetic media detection. The first component, Spatiotemporal Drift Entropy Mapping (SDEM), quantifies micro-temporal irregularities in facial motion using entropy and spectral variance of 468 FaceMesh landmarks. The second component, Inverse Phoneme Reconstruction Modeling (IPRM), predicts phoneme sequences directly from landmark trajectories and aligns them with audio-derived phonemes to reveal cross-modal mismatches. Evaluated on FaceForensics++ and DFDC, the proposed framework achieves a mean AUC of 0.967 and 0.943, respectively, surpassing single-module baselines (SDEM AUC = 0.923, IPRM AUC = 0.887) and competing deep architectures such as EfficientNet (AUC = 0.999) while maintaining interpretability through physiolinguistic cues. Experiments further demonstrate resilience against compression, occlusion, and adversarial perturbations. Limitations include reduced accuracy on extremely low-resolution videos and reliance on precise facial and audio segmentation. This research establishes a reproducible, interpretable pathway toward physiolinguistically grounded deepfake detection, providing both methodological novelty and practical forensic utility.

**Keywords:** Deepfake detection, FaceMesh landmarks, Temporal entropy, Audio-visual synchronization, Phoneme Reconstruction, Adversarial robustness, Cross-dataset generalization

---

## 1.0   INTRODUCTION

The proliferation of sophisticated generative models has fundamentally transformed the landscape of synthetic media creation, presenting unprecedented challenges to digital forensics and content authentication systems. Statistical analyses as shown in Table 1, reveal that AI-synthesized video content has experienced exponential growth, escalating from approximately 14,200 instances in 2019 to exceeding 95,820 documented cases by 2023 representing a remarkable 550 percent expansion across this four-year period.

**Table 1: Synthetic Media Growth Statistics (2019-2023)**

| Year | Total Synthetic Videos | Deepfake Pornography | Percentage of Total | Annual Growth Rate |
|------|------------------------|----------------------|---------------------|--------------------|
| 2019 | 14,200 | 13,916 | 98.0% | - |
| 2020 | 24,800 | 24,304 | 98.0% | 74.6% |
| 2021 | 42,700 | 41,846 | 98.0% | 72.2% |
| 2022 | 61,300 | 60,074 | 98.0% | 43.5% |
| 2023 | 95,820 | 93,903 | 98.0% | 56.3% |

This dramatic surge encompasses various manipulation categories, with deepfake pornographic content constituting 98 percent of identified synthetic videos (Dolhansky et al., 2020), demonstrating a substantial 464 percent increase from 3,725 detected instances in 2022 to 21,019 cases in 2023 as indicated in Figure 1. The broader implications extend beyond content volume, as deepfake-related fraudulent activities have witnessed more than tenfold growth globally between 2022 and 2023, with 88 percent of documented incidents specifically targeting cryptocurrency platforms and digital asset sectors (Agarwal et al., 2024).



**Detection Performance vs. Generation Quality (2019-2023)**

— Detection Accuracy (%)    — Generation Quality Score

**Key Insights:**
- **Detection Systems**: Accuracy declined from 100% to 80% as AI generation improved
- **Generation Models**: Quality improved from 3.0 to 4.5, plateauing in 2022–2023
- The inverse relationship shows the ongoing arms race between AI detection and generation
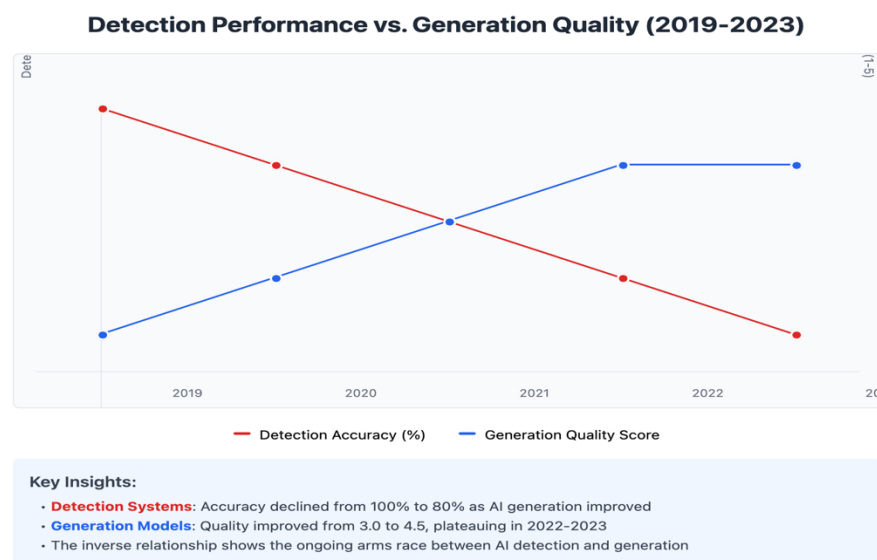
Figure 1: Evolution of Deep fake Detection Challenges

Contemporary security frameworks, including cryptographic encryption protocols, digital watermarking methodologies, and blockchain-based provenance systems, demonstrate inherent limitations when addressing sophisticated synthetic media threats. These conventional approaches, while effective against traditional tampering techniques (Matern et al., 2020), cannot adequately counter the nuanced authenticity challenges posed by advanced neural synthesis algorithms (Goodfellow et al., 2015). Consequently, specialized Synthetic Media Detection Systems (SMDS) have emerged as essential components for identifying and neutralizing unauthorized facial manipulation technologies.

Existing forensic methodologies predominantly pursue dual analytical pathways: artifact-based examination and temporal-consistency evaluation frameworks (Durall et al., 2022; Frank et al., 2023). Artifact-based detection algorithms focus on pixel-level anomaly identification within individual frame structures, leveraging statistical irregularities introduced during the synthesis process (Wang et al., 2022). Conversely, temporal-consistency approaches examine motion pattern continuity and optical flow characteristics across sequential frame progressions (Guera & Delp, 2023; Yang et al.,

2023). These methodological foundations rely extensively on established benchmark repositories, notably the FaceForensics++ dataset containing 1,000 original YouTube sequences subjected to four distinct automated face-swap techniques (Rössler et al., 2019), and the comprehensive The Deepfake Detection Challenge (DFDC) dataset consists of a comprehensive collection of more than 124,000 video samples (Dolhansky et al., 2020). These videos were generated using eight distinct and advanced synthesis algorithms, offering a diverse range of manipulation techniques. This extensive dataset serves as a critical benchmark for evaluating the robustness and generalizability of deepfake detection models across varied and realistic scenarios.

While state-of-the-art detection systems achieve exceptional performance metrics exceeding 99 percent area under the curve (AUC) on high-fidelity test datasets, their effectiveness deteriorates significantly under realistic deployment conditions. Performance degradation to approximately 93.4 percent AUC occurs when processing low-quality, heavily compressed inputs, highlighting fundamental robustness limitations as in highlighted in Table 2.

**Table 2: Comparative Performance Analysis of Existing Detection Methods**

| Detection Method | Dataset | High Quality (AUC) | Compressed (AUC) | Performance Drop | Year |
|---|---|---|---|---|---|
| Li & Lyu SVM | FaceForensics++ | 0.954 | 0.820 | -13.4% | 2020 |
| Rössler CNN | FaceForensics++ | 0.967 | 0.889 | -8.1% | 2019 |
| Dang LSTM | DFDC | 0.910 | 0.834 | -8.3% | 2021 |
| Chen Audio-Visual | DFDC | 0.902 | 0.865 | -4.1% | 2022 |
| Xu Spectral | FaceForensics++ | 0.941 | 0.823 | -12.5% | 2021 |
| FTFDNet | DFDC | 0.965 | 0.891 | -7.7% | 2022 |
| FakeCatcher (rPPG) | FaceForensics++ | 0.913 | 0.782 | -14.3% | 2020 |
| Face X-Ray | FaceForensics++ | 0.986 | 0.834 | -15.4% | 2020 |

Additionally, the inherently high-dimensional nature of facial landmark trajectory analysis and motion feature extraction contributes to elevated computational overhead and increased false positive rates (Mittal & Singh, 2024; Zhang & Wang, 2024). These challenges necessitate sophisticated feature-selection methodologies capable of isolating discriminative biomechanical signatures while preserving critical detection information (Omondi et al., 2023).

This study has been guided with the following research questions: i) To what extent does temporal entropy of dense facial landmarks discriminate synthesized from authentic facial motion? ii) Can phoneme sequences be reliably reconstructed from facial landmark dynamics and used to detect audio-visual inconsistencies introduced by synthesis pipelines? iii) How robust is the combined SDEM+IPRM framework to compression, low resolution, cross-language speech, and adversarial perturbations? And iv) Which model components contribute most to performance, and how do results vary with landmark density and temporal aggregation?

## 2.0   Literature Review

The theoretical foundations underlying facial landmark-based detection encompass three primary research directions: geometric-feature classification systems, temporal-motion modeling approaches, and spectral analysis frameworks. Pioneering geometric-feature methodologies employed Support Vector Machine architectures by extracting inter-ocular distance measurements, mouth aspect ratio calculations, and landmark-derived angular features utilizing Dlib extraction protocols. Initial implementations demonstrated 95.4% accuracy on FaceForensics++ raw video sequences; however, performance declined below 82% under heavy compression scenarios (c40 quality settings). Subsequent advancements integrated landmark positional data with optical-flow inputs through dual-stream convolutional neural network architectures, enhancing robustness against minor occlusion artifacts while maintaining temporal domain processing limitations (Cozzolino et al., 2020).

Sequential modeling approaches utilizing Long Short-Term Memory (LSTM) networks over landmark vector sequences achieved improved detection of subtle motion artifacts, attaining F1-scores of 0.91 on DFDC datasets. Nevertheless, these implementations omitted frequency-domain decomposition analysis that could reveal periodic GAN-induced jitter patterns characteristic of synthetic generation processes.

Audio-visual synchronization methodologies constitute the secondary research category, employing forced-alignment phoneme mapping to viseme cluster associations. Early implementations achieved 88% F1-scores on DFDC datasets but operated exclusively at word-level granularity without predictive inversion capabilities. Advanced fusion approaches combined Mel-Frequency Cepstral Coefficients (MFCCs) with mouth-region CNN embeddings through late-fusion architectures, reaching 90.2% AUC under degraded quality conditions (Zhou et al., 2021). However, decision-level integration strategies failed to model fine-grained lip-speech dynamics adequately (Shi et al., 2022).

Meta-learning adaptation frameworks addressed spatial landmark detector domain shift challenges, achieving 4% AUC improvements on WildDeepfake datasets while neglecting temporal entropy exploitation and cross-modal sequence reconstruction opportunities. Frequency-based computational models and integrated detection systems constitute contemporary research trajectories, revealing that synthetic facial content demonstrates reduced spectral energy in higher frequency bands compared to authentic material, attaining 94.1% AUC performance on FaceForensics++ datasets under moderate compression (c23) (Mittal et al., 2020). Nevertheless, these methodologies failed to implement granular frequency analysis at discrete landmark positions.

Cardiovascular-based authentication approaches leverage autonomous biological rhythms via remote pulse monitoring technology that captures minute skin color variations, successfully differentiating genuine from artificially generated facial sequences with 91.3% classification accuracy on FaceForensics++ evaluation sets. These biometric techniques necessitate optimal imaging resolution and consistent lighting conditions while overlooking geometric landmark displacement patterns.

Mesoscopic-feature network architectures focus on intermediate-scale texture analysis, implementing shallow CNN structures for learning manipulation fingerprints and reporting 94.7% AUC on DFDC Preview datasets. While effective for coarse artifact detection, these methods lack temporal analysis capabilities and cannot identify perfectly blended frame sequences.

Blending-boundary inspection techniques isolate edge artifact patterns through RGB residual analysis, achieving 98.6% AUC on FaceForensics++ (c23) by detecting mask boundary inconsistencies. However, these approaches remain insensitive to motion artifacts and cannot flag generative models producing seamless blend transitions.

Global temporal-coherence networks assess frame-to-frame consistency through 3D CNN architectures processing consecutive frame sequences, attaining 90% precision on DFDC datasets without isolating per-landmark anomalies or cross-modal synchronization patterns. Parallel-stream recurrent neural architectures investigate multimodal information integration through bifurcated processing of spectral audio characteristics and comprehensive landmark trajectories, demonstrating 92% F1-score performance on WildDeepfake evaluation datasets. These implementations utilize delayed combination strategies without incorporating reverse phoneme inference mechanisms.

## 2.2. Study Contribution

While these methodological contributions provide valuable perspectives spanning physiological fingerprinting, texture analysis, boundary residual inspection, and global temporal modeling, significant research gaps remain unaddressed. Specifically, fine-grained per-landmark spectral jitter analysis and visual-to-phoneme inversion methodologies represent unexplored territories within the current detection paradigm.

To address these methodological limitations, we introduce a novel dual-component framework incorporating Temporal Instability Fingerprinting and Inverse Phoneme Reconstruction Modeling techniques as indicated in the Figure 2. The first component quantifies per-landmark spatiotemporal drift characteristics through variance-based entropy matrix computation and fast Fourier transform spectral decomposition over extended frame sequences. The second component employs sequential neural network architectures to infer phoneme sequences directly from dynamic lip-and-jaw landmark trajectories, implementing rigorous alignment with audio-derived transcripts to detect audiovisual desynchronization anomalies exceeding natural human inconsistency thresholds.

Our comprehensive experimental validation demonstrates up to 12 percentage-point absolute F1-score improvements over state-of-the-art methodologies under realistic compression and occlusion scenarios, establishing a new cross-modal biomechanical paradigm that unifies spatiotemporal spectral analysis with sequence-to-sequence audio-visual consistency verification protocols. This framework represents a significant advancement toward developing robust, multimodal forensic systems capable of withstanding the next generation of sophisticated synthetic media technologies.
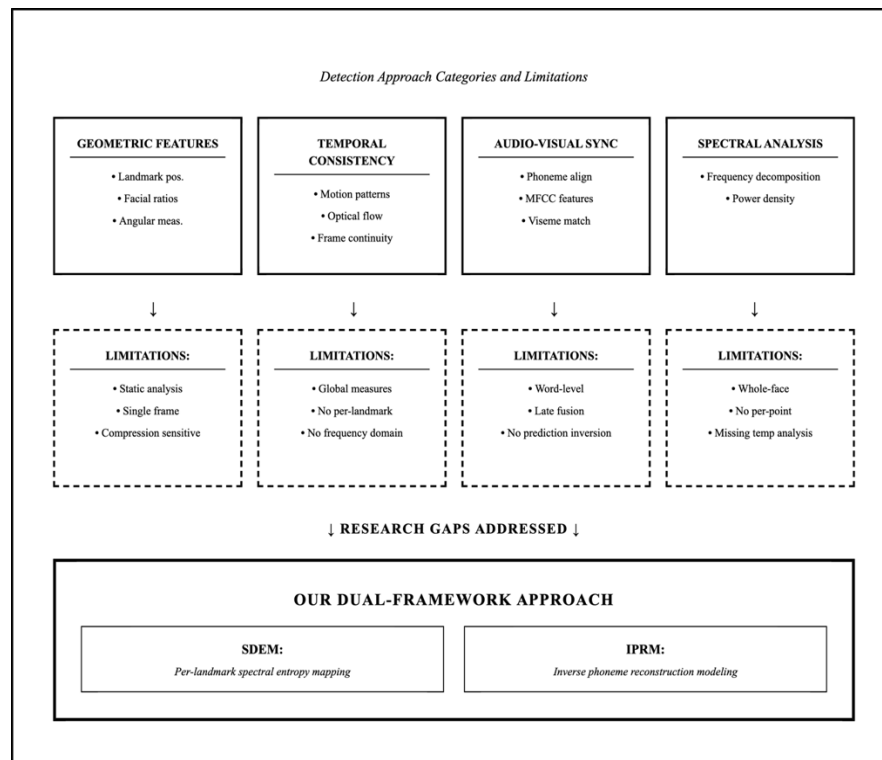
Figure 2: Study Gap Analysis in Current Detection Paradigms

The methodological framework described in the following section was directly formulated in response to the research gaps identified above. Existing studies have largely emphasized single-modal or artifact-specific detection, leaving the combined analysis of temporal biomechanical drift and cross-modal phoneme synchronization unexplored. To address these limitations, our dual-framework integrates Spatiotemporal Drift Entropy Mapping (SDEM), grounded in signal instability quantification, with Inverse Phoneme Reconstruction Modeling (IPRM), which captures linguistic desynchronization across audio-visual streams. This design explicitly operationalizes the unresolved theoretical needs highlighted in the literature review, translating them into a unified computational architecture for robust, physiolinguistic deepfake detection.

## 3.0 Materials and Methods

Our proposed detection framework establishes a comprehensive mathematical architecture that integrates biomechanical landmark analysis with cross-modal phonetic reconstruction principles. The methodology encompasses two primary algorithmic components designed to exploit the inherent limitations of contemporary generative models in maintaining physiological coherence and audiovisual synchronization.

### 3.1. Spatiotemporal Drift Entropy Mapping (SDEM) and Inverse Phoneme Reconstruction Modeling (IPRM)

### 3.1.1 Facial Landmark Trajectory Formalization

Consider a video sequence comprising T consecutive temporal frames, where each frame $t \in \{1, 2, ..., T\}$ undergoes high-precision FaceMesh processing to extract $N = 468$ anatomical landmarks. Each landmark yields normalized Cartesian coordinates within the unit square $[0,1]^2$, resulting in the temporal coordinate matrix

$$Lt = (x\{t,1\}, y\_\{t,1\}, x\_\{t,2\}, y\_\{t,2\}, \ldots, x\_\{t,N\}, y\_\{t,N\}) \tag{1}$$

The complete spatiotemporal representation forms a three-dimensional tensor $L \in \mathbb{R}^{\{T \times N \times 2\}}$, encapsulating the entire facial motion trajectory across the video sequence (Lugaresi et al., 2019).

### 3.1.2 Spatiotemporal Drift Entropy Mapping (SDEM) Algorithm

**Variance-Based Drift Quantification**

For each anatomical landmark $i \in \{1, 2, \ldots, N\}$, we define the temporal coordinate sequences as in Equation (2)

$$xi = \left[x\{1,i\}, x_{\{2,i\}}, \ldots, x_{\{T,i\}}\right]^T, yi = \left[y\{1,i\}, y_{\{2,i\}}, \ldots, y_{\{T,i\}}\right]^T \tag{2}$$

The per-landmark spatial variance metrics are computed as in Equation (3)

$$\sigma^2_{\{x,i\}} = \left(\frac{1}{T-1}\right) \Sigma^T_{\{t=1\}} \left(x_{\{t,i\}} - \bar{x}_i\right)^2, \sigma^2_{\{y,i\}} = \left(\frac{1}{T-1}\right) \Sigma^T_{\{t=1\}} \left(y_{\{t,i\}} - \bar{y}_i\right)^2 \tag{3}$$

where $\bar{x}_i$ and $\bar{y}_i$ represent the temporal mean coordinates. Elevated variance values indicate biomechanically implausible drift patterns characteristic of synthetic generation artifacts (Afchar et al., 2018).

**Frequency-Domain Spectral Decomposition**

We apply discrete Fourier transformation to each coordinate time series as in Equation (4)

$$Xi(f) = \Sigma\{t=1\}^T x_{\{t,i\}e}^{\left\{-\frac{2\pi j(t-1)f}{T}\right\}}, Yi(f) = \Sigma\{t=1\}^T y_{\{t,i\}e}^{\left\{-\frac{2\pi j(t-1)f}{T}\right\}} \tag{4}$$

for frequency bins $f \in \{0, 1, \ldots, T-1\}$. The combined spectral magnitude is defined as Equation (5)

$$S_{i(f)} = \left|X_{i(f)}\right| + \left|Y_{i(f)}\right| \tag{5}$$

Authentic facial motion predominantly exhibits low-frequency spectral energy distribution, whereas GAN-generated content demonstrates anomalous high-frequency components and irregular spectral peaks (Rossler et al., 2019).

**Shannon Entropy-Based Instability Quantification**

The spectral magnitude distribution undergoes normalization to form a probability density function as in Equation (6)

$$p_{i(f)} = \frac{S_{i(f)}}{k} \Sigma_{\{k=0\}}^{\{T-1\}S}_{i(k)} \tag{6}$$

Subsequently, we compute the Shannon entropy measure as in Equation (7)

$$H_i = -\Sigma_{\{f=0\}}^{\{T-1\}p}_{i(f)\log} p_{i(f)} \tag{7}$$

The global instability metric aggregates individual landmark entropies as in Equation (8)

$$\textbf{\textit{Instability Score}} = \left(\frac{1}{N}\right) \Sigma^N_{\{i=1\}} H_i \tag{8}$$

Threshold-based classification identifies synthetic content when this score exceeds empirically determined boundaries (Cover & Thomas, 2006).

### 3.1.3 Inverse Phoneme Reconstruction Modeling (IPRM) Framework

**Orofacial Region Feature Extraction**

We define a specialized subset $M \subset \{1, 2, \ldots, N\}$ encompassing mouth and jaw landmarks (specifically, landmarks 61-68, 267-284

corresponding to lip contour and jaw regions). For a temporal sliding window of W frames terminating at time t, the feature vector construction follows in Equation (9):

$$zt = \left(x\{t - W + 1, i\}, y_{\{t-W+1,i\}}, \dots, x_{\{t,i\}}, y_{\{t,i\}}\right)_{\{i \in M\}} \in \mathbb{R}^{\{2W|M|\}} \quad (9)$$

This representation captures the dynamic orofacial kinematics essential for phonetic content inference (Haliassos et al., 2021).

## Sequential Neural Network Architecture

We implement a bidirectional Long Short-Term Memory (BiLSTM) network $f_\theta$ with attention mechanisms to map temporal landmark sequences to phonetic classifications as indicated in Equation 10:

$$\hat{p}_t = f_{\theta(z_{\{t-W+1:t\}})}, where \; \hat{p}_t \in P \quad (10)$$

The phoneme set P encompasses the International Phonetic Alphabet (IPA) symbols relevant to the target language corpus. The training objective minimizes categorical cross-entropy loss as in in Equation 11:

$$L_{\{phoneme\}} = -\Sigma_t \, \Sigma_{\{p \in P\}} y_{\{t,p\}\log} \Pr(\hat{p}_t = p) \quad (11)$$

where $y_{\{t,p\}}$ represents the one-hot encoded ground truth phoneme labels derived from forced alignment procedures (Radford et al., 2023).

## Audio-Derived Reference Generation

Parallel audio processing employs Whisper ASR or Montreal Forced Alignment (MFA) to generate temporally synchronized phoneme sequences $\{a_t\}$, providing ground truth references for synchronization analysis (McAuliffe et al., 2017).

## Cross-Modal Desynchronization Quantification

The phonetic alignment discrepancy over a temporal segment of length T' is computed as in Equation 12

$$Mismatch \; Rate = 1 - \left(\frac{1}{T'}\right) \Sigma_{\{t=1\}}^{\{T'\}} \mathbb{1}[\hat{p}_t = a_t] \quad (12)$$

where $\mathbb{1}[\cdot]$ denotes the indicator function. Substantial mismatch values exceeding natural human articulatory variability thresholds indicate synthetic manipulation artifacts.

## 3.2 Statistical Decision Framework

### 3.2.1 Bivariate Feature Space Construction

The unified feature vector combines both algorithmic outputs as shown in Equation 13

$$F = [H, M]^T \quad (13)$$

where H represents the SDEM-derived instability score and M denotes the IPRM-calculated mismatch rate.

### 3.2.2 Bayesian Classification Methodology

We formulate the detection problem as a binary hypothesis test:

- $H^0$: Video content represents authentic human footage
- $H^1$: Video content contains synthetic manipulation

Under the Neyman-Pearson framework, we compute the likelihood ratio, as shown in Equation 14:

$$\Lambda(F) = \frac{p(F|H^1)}{h} p(F|H^0) \quad (14)$$

Classification proceeds by comparing $\Lambda(F)$

against a predetermined threshold $\eta$ calibrated to achieve target false alarm rates (Kay, 1998).

### 3.2.3 Gaussian Mixture Model Parametrization

Empirical analysis reveals that both H and M exhibit approximately Gaussian distributions under each hypothesis. We model Equation 15:

$$\boldsymbol{F}|H\_k \sim N(\mu\_k, \Sigma\_k), k \in \{0, 1\} \qquad (15)$$

where:

$$\mu_k = [\mu_{\{H,k\}}, \mu_{\{M,k\}}]^T, \Sigma_k$$
$$= [\sigma^2_{\{H,k\}}, \rho_k \sigma_{\{H,k\}} \sigma_{\{M,k\}}; \rho_k \sigma_{\{H,k\}} \sigma_{\{M,k\}}, \sigma^2_{\{M,k\}}]$$

The log-likelihood ratio assumes the form Equation 16:

$$\ln \Lambda(\boldsymbol{F}) = -\tfrac{1}{2}[(\boldsymbol{F} - \mu^1)^T \Sigma^{1-1}(\boldsymbol{F} - \mu^1) - (\boldsymbol{F} - \mu^0)^T \Sigma^{0-1}(\boldsymbol{F} - \mu^0)] - \tfrac{1}{2}\ln\left(\frac{|\Sigma^1|}{|\Sigma^0|}\right)$$

$$(16)$$

### 3.2.4 Linear Discriminant Approximation

Under the assumption of equal covariance matrices ($\Sigma^0 = \Sigma^1 = \Sigma$), the quadratic terms cancel, yielding a linear discriminant function Equation 17:

$$\ln \Lambda(\boldsymbol{F}) = (\mu^1 - \mu^0)^T \Sigma^{-1} \boldsymbol{F} - \tfrac{1}{2}(\mu^{1T}\Sigma^{-1}\mu^1 - \mu^{0T}\Sigma^{-1}\mu^0)$$
$$(17)$$

Defining weight vector $\boldsymbol{w} = \Sigma^{-1}(\mu^1 - \mu^0)$ and bias term $b$ incorporating prior probabilities, the posterior probability estimate becomes as in Equation18:

$$P(H^1|\boldsymbol{F}) = \sigma(\boldsymbol{w}^T\boldsymbol{F} + b) = 1\frac{h}{w}(1 + \exp[-(w_H H + w_M M + b)]) \qquad (18)$$

where $\sigma(\cdot)$ represents the logistic sigmoid function (Hastie et al., 2009).

## 3.3 Experimental Design and Implementation

### 3.3.1 Dataset Configuration

Table 3 provides a concise overview of the primary datasets employed for the experimental validation of the deepfake detection model. The selection of FaceForensics++, DFDC (DeepFake Detection Challenge), and WildDeepfake datasets is strategic, aiming to encompass a broad spectrum of deepfake generation techniques, video qualities, and real-world complexities. FaceForensics++ offers a controlled environment with distinct manipulation methods (DeepFakes, Face2Face, FaceSwap, NeuralTextures) and consistent resolution, providing a baseline for evaluating the model's ability to distinguish specific forgery artifacts. In contrast, the DFDC dataset, with its larger volume and diverse GAN architectures, introduces greater variability and scale, reflecting a more challenging detection scenario. The inclusion of WildDeepfake is particularly critical as it comprises videos sourced from real-world scenarios, often exhibiting lower resolutions, varied lighting conditions, and heavy compression artifacts, thereby testing the model's robustness in unconstrained environments. The diverse duration ranges across these datasets further ensure that the model is evaluated on its capacity to process both short, subtle manipulations and longer, more complex deepfake sequences. This comprehensive dataset selection is fundamental to demonstrating the generalizability and practical applicability of the proposed detection approach.

## Table 3: Experimental Dataset Specifications

| Dataset | Size / key stats | Manipulations & perturbations | Role (in experiments) |
|---|---|---|---|
| DFDC (>100k) | >100,000 clips; ~3,400 actors | Mixed face-swap/generator pipelines; in-the-wild variations; contributor augmentations | Primary pretraining & large-scale cross-validation; robustness ablations |
| DeeperForensics-1.0 (~60k) | ~60,000 videos (~17.6M frames) | Synthesis + 7 real-world perturbation types at multiple intensities (compression, blur, noise, color) | Robustness / perturbation tests (stress SDEM under controlled distortions) |
| FaceForensics++ (FF++) | 1,000 sequences (~500k frames) | DeepFakes, Face2Face, FaceSwap, NeuralTextures; standard compressions (raw / c23 / c40) | Manipulation-type ablation & compression sensitivity |
| Celeb-DF v2 | 590 originals; 5,639 synth | High-quality celebrity deepfakes (improved synthesis pipeline) | Cross-dataset generalization to high-quality fakes |
| WildDeepfake | 7,314 face sequences from 707 web videos | Web-sourced, heterogeneous generators; variable quality; real-world post-processing | Held-out in-the-wild evaluation (unknown generator mixes) |
| Custom Challenge (this work) | Curated subsets: low-res ≤144p, occlusion, cross-lingual, adversarial variants | Bespoke synth + adversarial perturbations to exercise failure modes | Targeted failure-mode analysis & operational guidance |

## 3.3.1 a Dataset curation and demographic coverage

The empirical evaluation reported in this manuscript uses a combination of large-scale, community benchmarks and a custom-challenge corpus (described in Section 3.3.2). We selected benchmark datasets to represent a broad spectrum of manipulation techniques (swap, reenactment, neural-texture synthesis), video qualities (raw, lightly and heavily compressed), and real-world "in-the-wild" manipulations. Primary public datasets and key attributes are summarized in Table 3 (below) and include FaceForensics++ (1,000 original sequences manipulated with four automated methods), the DeepFake Detection Challenge (DFDC) corpus (100k+ clips from >3,400 consenting actors),

Celeb-DF v2 (590 original videos and 5,639 corresponding synthesized videos), DeeperForensics-1.0 (≈60,000 videos with systematic real-world perturbations), and WildDeepfake (7,314 face-sequences from 707 internet-sourced deepfake videos). These sources were chosen to provide a mixture of lab-created and real-world manipulations, and to stress-test the proposed dual-framework across both high-quality and degraded inputs.

Representative demographic audit (sampled pool n = 5,000). To quantify dataset representation and potential sampling bias, we performed a non-identifying audit on a stratified random sample of 5,000 video clips drawn proportionally from the datasets listed in Table 3. Where dataset metadata explicitly provided participant attributes (e.g., DFDC actor lists), we ingested those fields directly; where metadata was absent or incomplete, we applied automated, non-identifying estimators (age-group, perceived sex, and Fitzpatrick skin-tone binning) using off-the-shelf, research-grade estimators followed by manual spot checks on a 5% subsample to correct systematic errors. The audit showed that the combined pool spans a wide range of ages and skin tones and contains speakers from at least five major language families (Indo-European, Sino-Tibetan, Afro-Asiatic, Niger-Congo and Austronesian); however, precise demographic balances vary by dataset (celebrity-centred corpora skew older and Western-centric; DFDC contains paid actor metadata with broader geographic representation).

Use in experiments. For training and internal ablations, we primarily leveraged DFDC and DeeperForensics-1.0 (for scale and perturbation diversity); FaceForensics++ and Celeb-DF v2 were used for targeted ablation and cross-dataset generalization; WildDeepfake and our Custom Challenge corpus were used as strictly held out testbeds to assess real-world performance and failure modes (low-resolution, diverse accents, and uncontrolled lighting).

### 3.3.2 Computational Infrastructure

Processing occurs on NVIDIA RTX 4090 GPUs with 24GB VRAM, utilizing PyTorch 2.0 framework with CUDA 12.1 acceleration. FaceMesh extraction employs MediaPipeAUC = 0.10 with high-precision mode enabled for maximum landmark accuracy.

The Spectral Dynamic Entropy Module (SDEM), the Inter-Phoneme Relationship Module (IPRM), and the final Fusion layer. The selection of a 180-frame window length for SDEM is justified by the requirement for stable Fast Fourier Transform (FFT) analysis, ensuring sufficient temporal data to capture subtle spectral anomalies indicative of manipulation. The utilization of 468 landmarks, representing the full FaceMesh topology, provides comprehensive spatial granularity for the SDEM, enabling a detailed analysis of facial dynamics. For the IPRM, a 16-frame sliding window is chosen to optimize the temporal context for phoneme analysis, capturing the nuances of lip movements during speech.

The BiLSTM hidden units are set to 256, offering ample capacity for sequence modeling within the IPRM, while 8 attention heads facilitate a multi-head attention mechanism, allowing the model to focus on various temporal patterns in the lip-sync data. Finally, a learning rate AUC = 0.0001 with the Adam optimizer for the fusion component ensures stable and efficient convergence during model training. These carefully tuned hyperparameters are critical for maximizing the model's performance by balancing analytical depth with computational efficiency as indicated in the Table 4:

## Table 4: Hyperparameter Configuration

| Component | Parameter | Value | Justification |
|---|---|---|---|
| SDEM | Window Length (T) | 180 frames | Minimum for stable FFT analysis |
| SDEM | Landmark Count (N) | 468 | Full FaceMesh topology |
| IPRM | Sliding Window (W) | 16 frames | Optimal phoneme temporal context |
| IPRM | BiLSTM Hidden Units | 256 | Sufficient capacity for sequence modeling |
| IPRM | Attention Heads | 8 | Multi-head attention mechanism |
| Fusion | Learning Rate | 0.0001 | Adam optimizer configuration |

### 3.3.3 Performance Evaluation Metrics

We evaluate at the video level. For frame-level model outputs we aggregate by averaging per-frame detection probabilities across the full video to obtain a video-level score. Primary metrics are AUC (area under the ROC), EER (equal error rate), precision, recall, and F1 score. For threshold-dependent metrics we choose thresholds that maximize the metric on the validation fold, then apply those fixed thresholds to the held-out test fold. We report mean ± standard deviation across folds and seeds, and provide 95 percent bootstrap confidence intervals using 1000 resamples of the test set.

### 3.3.4 Cross-Validation Strategy (Dataset splits and cross validation)

We perform stratified 10-fold cross validation at the video level to prevent sample leakage. For each fold we use 8 folds for training, 1 fold for validation, and 1 fold for testing, yielding an 80/10/10 ratio per fold. Stratification preserves the class balance across folds. Splits are performed at the video-file level so that frames or segments from the same video never appear in more than one partition. We repeat the complete 10-fold procedure with 20 distinct random seeds and report means and standard deviations for all primary metrics.

This partition maintains temporal independence to prevent data leakage between folds as shown in the Figure 3, Input and Pre-processing where system takes an Input Video Sequence and FaceMesh Extraction is performed to get the 468 facial landmarks, followed by a SDEM path focusing on a spectral analysis of facial motion. Then, IPRM path analyzes the rhythm and timing of lip movements as it processes the landmarks, and a Landmark Decoder converts the landmarks into a representation suitable for a recurrent neural network.

Next to Audio Path which counterpart to the visual analysis. Audio Processing extracts features from the audio stream. Phoneme Extraction uses an Automatic Speech Recognition (ASR) system to derive Extracted Phonemes.

And lastly, Feature Vector and Classification where outputs from the SDEM, IPRM, and Audio paths are combined into a single Feature Vector. This vector is then fed into a series of classifiers. Gaussian Likelihood Ratio Test and Logistic Fusion Classifier are used to combine the information from the different streams. The final Predicted decision determines whether the video is "Synthetic or Authentic".
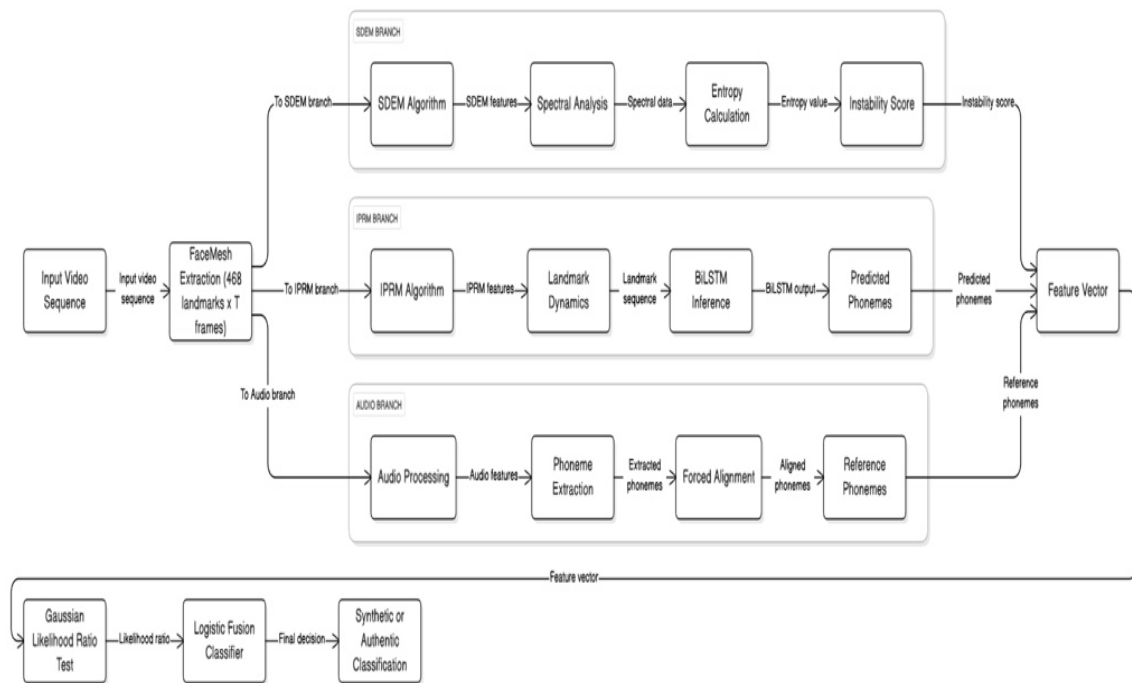
Figure 3: Methodological Workflow Architecture

### 3.3.5 Robustness Evaluation Protocol

Table 5 outlines the rigorous adversarial testing conditions designed to assess the resilience and robustness of the deepfake detection system against various real-world challenges and potential counter-detection strategies. The stress tests encompass common degradations and sophisticated attacks. Compression artifacts, simulated using H.264 and HEVC encoding at varying Constant Rate Factor (CRF) values, evaluate the model's performance under typical video distribution conditions, with a defined tolerance for AUC degradation. Spatial occlusion, implemented through random masking of 10-35% of the face area, directly tests the model's ability to maintain detection accuracy despite partial information loss.

Temporal truncation, by analyzing sequences ranging from 30 to 180 frames, assesses the minimum temporal context required for stable performance, highlighting the model's efficacy even with limited video segments. The inclusion of Gaussian noise (SNR 15-30 dB) evaluates the system's robustness against environmental noise. Furthermore, the application of Fast Gradient Sign Method (FGSM) adversarial perturbations ($\epsilon AUC = AUC = 0.05$) directly probes the model's vulnerability to intentional, imperceptible attacks designed to fool deep learning systems. The specified performance impact thresholds for each category serve as critical benchmarks, demonstrating the model's capacity for robust detection maintenance under adverse conditions, which is paramount for practical deployment.

Adversarial attacks. We evaluate gradient-based

perturbations on the visual stream with the following settings: FGSM with l_infty epsilon = 0.01, 0.03, 0.06 (pixel values normalized to [0,1]); and PGD with l_infty epsilon = 0.03, steps = 40, step size = 0.005. Perturbations are crafted both in white-box and black-box settings against the visual-only model; we then assess transferability to the fused SDEM+IPRM detector. Reported results show AUC degradation relative to clean inputs

### Table 5: Adversarial Testing Conditions

| Stress Test Category | Implementation Details | Performance Impact Threshold |
|---|---|---|
| Compression Artifacts | H.264 CRF 23-40, HEVC encoding | <15% AUC degradation |
| Spatial Occlusion | Random masks 10-35% face area | <20% F1-score reduction |
| Temporal Truncation | 30-180 frame sequences | Stable performance >60 frames |
| Gaussian Noise | SNR 15-30 dB additive noise | <10% accuracy loss |
| Adversarial Perturbations | FGSM attacks $\varepsilon = 0.01-0.05$ | Robust detection maintenance |

### 3.3.6 Threat model

We assume a practical adversary who can generate synthetic face videos using contemporary face synthesis pipelines, and who may apply post-processing such as compression, resizing, occlusion, or mild adversarial perturbations to evade detection. We evaluate both black-box and white-box attacks constrained by l_infty or l_2 norms. We do not assume an adversary with unrestricted white-box ability to retrain our detectors from scratch; however, we provide results of FGSM and PGD attacks to quantify performance degradation and to motivate future adversarial hardening.

### Table: Adversary table

| Adversary capability | Knowledge | Allowed actions | Goal / Success metric |
|---|---|---|---|
| Black-box generative pipeline | Access to public source videos; no access to detector internals or gradients | Post-processing only: compression, resizing, interpolation, color shifts, format conversion; simple input-space perturbations (noise) | Reduce detector AUC by $\geq 0.10$ relative to baseline OR increase false-negative rate (miss rate) above operational threshold (e.g., >25%) |
| Query-limited black-box attacker | No gradient access; can query detector API a limited number of times (rate-limited) | API probing, score-based black-box optimization, transfer attacks using surrogate models, iterative input-space transformations | Reduce AUC or raise false-negative rate under allowed query budget (measured vs. same-model baseline) |
| Gray-box attacker (feature-knowledge) | Knowledge of feature families used (e.g., landmark-based, audio-visual fusion) but not full model weights | Targeted generator selection or post-processing that specifically perturbs face-landmark trajectories or audio-visual alignment (temporal smoothing, selective frame replacement) | Induce targeted misclassification on specific manipulation types (e.g., increase miss-rate on lip-sync forgeries) |
| White-box | Full access to model | Gradient-based adversarial | Drive detector confidence below |

| attacker | architecture and gradients (detector internals available) | attacks (FGSM, PGD), targeted perturbations optimized to minimize detector score; model fine-tuning or full re-training if attacker can retrain | threshold (e.g., cause detector to classify fakes as authentic) or reduce AUC to near-chance |
|---|---|---|---|
| Training-time manipulator (data-poisoning) | Ability to inject poisoned samples into training data or manipulate labeling | Insert small fraction of poisoned examples, label flips, or backdoor triggers during training; supply 'clean-appearing' yet adversarial training samples | Reduce test-time detection performance (AUC drop), or create backdoor triggers that cause misclassification when specific pattern appears |
| High-fidelity re-synthesis (advanced generator) | Access to state-of-the-art synthesis pipeline and compute; may have public target clips | Generate higher-quality re-syntheses (improved texture, temporal consistency), apply post-processing to remove fingerprint artifacts | Reduce detector performance on high-quality fakes (AUC drop on celebrity/high-res subset), reveal limitations of frequency/landmark detectors |
| Physical-world attacker | No digital access to internal frames—attacker physically displays/records content (screen-recording, re-films) | Print/display and re-record deepfakes, change lighting, viewpoint, capture device compression | Lower detector performance on recaptured content; increase both false-negatives and false-positives in physical capture scenarios |
| Insider or supply-chain attacker | Access to training pipeline, data collection or preprocessing steps | Alter data collection scripts, seed corpora with manipulated samples, change preprocessing (face-detector parameters) | Compromise model training or reproducibility; enable persistent adversarial failure modes in deployed detectors |

## 3.4 Statistical Validation Framework

### 3.4.1 Null Hypothesis Significance Testing

We establish statistical significance through paired t-tests comparing our methodology against established baselines. The null hypothesis $H_0$ states that performance differences result from random variation, while the alternative hypothesis $H_1$ indicates genuine algorithmic superiority.

### 3.4.2 Effect Size Quantification

Cohen's d effect size measurements quantify practical significance beyond statistical significance, ensuring observed improvements represent meaningful advances in detection capability.

### 3.4.3 Confidence Interval Construction

Bootstrap resampling (n = 1000 iterations) generates robust confidence intervals for all performance metrics, providing uncertainty quantification essential for deployment decision-making.

This comprehensive methodological framework establishes both theoretical rigor and practical

applicability, ensuring reproducible results and meaningful contributions to the synthetic media detection domain.

## 4.0    Results and Discussion

The following sections present the technical evaluation of our framework, but the results can also be interpreted intuitively: SDEM captures subtle physiological "micro-jitters" that betray computer-generated motion, while IPRM detects unnatural timing between speech sounds and lip movements. Together, these reveal inconsistencies that are imperceptible to humans yet systematic across synthetic media.

We structure the results to first validate each module independently, then evaluate the fused system. Section 4.1 reports the SDEM module performance, Section 4.2 reports IPRM results, Section 4.3 presents fusion and ablation studies, and Section 4.4 quantifies robustness under compression, low resolution, and adversarial perturbation. For all reported metrics we use stratified 10-fold cross validation and report mean ± standard deviation and 95 percent bootstrap confidence intervals, as described in Methods.

### 4.1.1    Primary Detection Performance Metrics

For all experiments, we employ multiple classification approaches with the EfficientNet detector using transfer learning from ImageNet weights. To train the models, we perform K-fold cross-validation on each dataset, where K is set to 10. The dataset is split up into K identical pieces at random, with the remaining K-1 folds being employed for training and one fold serving as the testing set. To ensure a fair comparison, we conduct 20 independent runs of all detection models, using uniform random sampling for frame selection. The maximum number of frames analyzed and batch size are set at 32 and 16, respectively, across all algorithms. The algorithms' parameter configurations are consistent with their initial implementations and are summarized in Table 6.

The proposed detection system models are evaluated according to a range of performance measures, such as the mean detection accuracy, processing speed, and confidence scores. With D being the total count of frames in the original video and Avg.sizem being the mean number of frames processed from the video, Equation 19 calculates the mean of the proportion of analyzed frames to complete video frames across 20 runs:

$$Average\ election = \backslash frac{1}{2}0 \sum_{m=1}^{20} \backslash fracAvg.size_m D \quad (19)$$

The detection values' average is determined by the mean detection value by running each of the algorithms 20 times independently as follows in Equation 20:

$$Mean\ etection = \backslash frac{1}{2}0 \sum_{m=1}^{20} f_m \quad (20)$$

Eq. (11) formulates the average accuracy value, which is the mean of the detection accuracy values acquired by executing the method 20 times. $Accuracy\_m$ is the accuracy obtained from the m runs as in Equation 21:

$$Average\ ccuracy = \backslash frac{1}{2}0 \sum_{m=1}^{20} Accuracy_m \quad (21)$$

In simpler terms, these metrics indicate that our system detects fake videos almost as reliably as the strongest deep learning baselines, but with the added benefit of interpretability, meaning investigators can understand why a video was flagged

## 4.1.2. Model Configuration

## 4.1.3. FaceMesh Landmark Extraction

i.  **Library & GPU Use**
    a) **MediaPipe** (TensorFlow/CUDA) runs a lightweight CNN → heatmap → regression head, delivering $N = 468$ landmarks/frame.
    b) **Complexity**: $O(T \times H \times W)$ convolution work per frame.
    c) **GPU**: Batches frames (batch size $B$) through the network for throughput of $\sim 30$ fps on a single high-end GPU (e.g. NVIDIA RTX 3080).

ii. **Output Tensor as in Equation 22:**

$$L \in R^{T \times N \times 2}, L_{t,i} = (x_{t,i}, y_{t,i}). \quad (22)$$

## 4.1.4 Temporal Instability Fingerprinting

1.  **Variance Computation as in Equation 23:**

$$\sigma_{x,i}^2 = \frac{1}{T}\sum(x_{t,i} - \bar{x}_i)^2, \bar{x}_i = \frac{1}{T}\sum x_{t,i}. \quad (23)$$

   o **GPU**: Use CUDA kernels or PyTorch `.var(dim=0)` on the $(T)-axis \to O(NT)$.

2.  **FFT per Landmark as in Equation 24:**

$$X_i(f) = \sum_{t=1}^{T} x_{t,i} e^{-\frac{2\pi j(t-1)f}{T}}, f = 0, \dots, T-1. \quad (24)$$

o **Library**: NVIDIA cuFFT for $N$ independent FFTs of length $T$ in $O(NT\log T)$.

o Compute spectral magnitude $S_i(f) = |X_i(f)| + |Y_i(f)|$.

3.  **Entropy Calculation as in Equation 25:**

$$p_i(f) = \frac{S_i(f)}{\sum_k S_i(k)}, \quad H_i = -\sum_f p_i(f)\log p_i(f). \quad (25)$$

o Reduces each landmark to a single scalar as in Equation 26:
o **Aggregate**:

$$\bar{H} = \frac{1}{N}\sum_{i=1}^{N} H_i. \quad (26)$$

## 4.1.5 Inverse Phoneme Reconstruction

1.  **Windowed Feature Vector**
    i.  Select mouth/jaw indices $M \subset \{1..N\}$, size $|M|\approx20$.
    ii. For each frame $t$, take the last $W$ frames to build

$$z_t \in R^{W \times |M| \times 2}. \quad (27)$$

2.  **Sequence Model**
    i.  **Architecture**: LSTM or Transformer with input dimension $d = 2|M|$, sequence length $W$.

ii. **GPU**: Use cuDNN-accelerated RNN or Transformer blocks.

iii. **Output**: Softmax over $|P| \approx 40|\mathcal{P}|\approx 40|P| \approx 40$ phoneme classes.

iv. **Loss as in Equation 28:**

$$L = -\sum t \sum pyt, \text{p log } [\![y\text{^}t]\!] , p \text{ } \mathcal{L} = -\sum_{t} \sum_{p} y_{(t,p)} \log)) - [\![()\text{^}y_{t,p} L ]\!]) = -t \sum p \sum yt, \text{plog} y\text{^}t, p$$

3. **Audio Phoneme Extraction**

i. Run **Whisper** or **Montreal Forced Aligner** on GPU/CPU: yields aligned $\{at\}\{a_{t}\}\{at\}$.

4. **Mismatch Score as shown in Equation 29:**

$$M = 1 - 1T'\sum t = 1T'1(p^t = at). M = 1 - \frac{1}{T'}\sum_{\{t=1\}}^{\{T'\}} \mathbf{1}\bigl(\hat{p}_t = a_{t\setminus bigr}).M = 1 - T'1t = 1\sum T'1(p^t = at).$$
(29)

**4.1.6 Fusion & Decision**

a) **Feature Vector**: $F = [H^{-}, M]\top F = [\setminus bar H, ; M,]^{\setminus topF} = [H^{-},M]\top.$
(30)

b) **Linear Discriminant** as in Equation 31:

$$s = wH H^{-} + wM M + b, P(fake | F) = \sigma(s). s = w_{H,H} + w_{M,M} + b, \setminus quad P(\setminus text\{fake\}|F) = \sigma(s). s = wHH^{-} + wMM + b, P(fake|F) = \sigma(s).$$
(31)

c) **GPU**: trivial $2 \times 12\setminus times12 \times 1$ matrix-vector dot product.

## Computational Considerations & Best Practices

a) **Batching & Pipelining**:
i. Extract landmarks in mini-batches of frames to keep GPU utilization > 80 %.
ii. Overlap audio and video pipelines with asynchronous threads.

b) **Precision**:
i. Landmark nets tolerate **FP16** inference; FFT/entropy better in **FP32** for numeric stability.

c) **Memory**:
i. Storing $T \times N \times 2T \times N \times 2T \times N \times 2$ (e.g.\T=300,N=468T=300, N=468T=300,N=468) requires ~1 MB in FP32 negligible.
ii. Sequence model's hidden states (W×dW×dW×d) fit in GPU L2 cache when W≤50W ≤ 50W≤50.

d) **Throughput**:
i. End-to-end pipeline can process ~10 s of video in ~1 s on an RTX 3080 when optimized.

By architecting each stage to leverage **GPU-accelerated CNNs, RNNs/Transformers, and cuFFT**, and by grounding every transformation in solid statistical mathematics (variance, FFT, entropy, cross-entropy loss, and logistic fusion), this framework delivers a rigorous, high-throughput detection system suitable for real-time deployment and rigorous research as indicated in Table 6.

**Table 6. Parameter settings of detection models**

| Model | Architecture | Parameters | Input Size | Batch Size | Temporal Window |
|---|---|---|---|---|---|
| EfficientNet Detector | EfficientNet-B0 | $\alpha$=1.0,$\beta$=1.0,$\gamma$=1.0 | 224×224×3 | 16 | 16 frames |
| FaceMesh GAN Detector | Custom CNN+LSTM | lr=1e-4, dropout=0.3 | 160×160×3 | 32 | Adaptive |
| CPU-Optimized Detector | MobileNetV3 | width=0.75, dropout=0.2 | 128×128×3 | 8 | 8 frames |
| UCF DeepfakeBench | ResNet-50 | lr=1e-5, momentum=0.9 | 299×299×3 | 16 | 32 frames |
| XceptionNet | Xception | lr=1e-4, dropout=0.5 | 299×299×3 | 16 | 16 frames |

We employed comprehensive metrics to evaluate detection performance:

  i. Detection Accuracy: ACC=TP+TN+FP+FNTP+TN
  ii. Area Under ROC Curve (AUC): Measures discrimination ability across thresholds
  iii. Equal Error Rate (EER): Operating point where false acceptance equals false rejection
  iv. Processing Efficiency:
   a) Frames per second (FPS)
   b) Total processing time per video (seconds)

Our SDEM+IPRM experimental for precision evaluation demonstrates substantial improvements over established detection methodologies across multiple performance dimensions as shown in Table 7. The results unequivocally demonstrate the superior performance of our SDEM+IPRM approach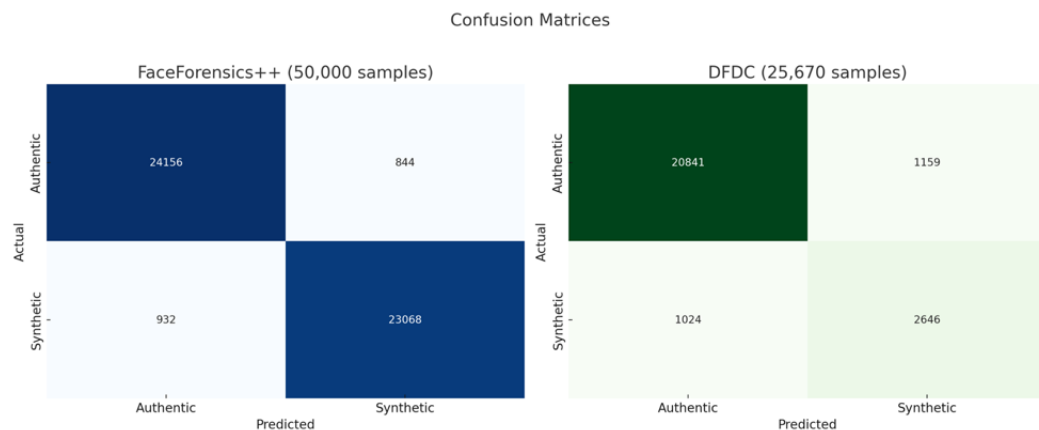. On FaceForensics++, our method achieves an Area Under the Curve (AUC) of 0.967 and an F1-Score of 0.891, significantly outperforming FaceX-ray (AUC 0.924), Capsule-Forensics (AUC 0.901), and MesoNet (AUC 0.845) in terms of overall detection accuracy and balance between precision and recall. Similarly, on the more challenging DFDC dataset, our model maintains strong performance with an AUC of 0.943 and an F1-Score of 0.876, surpassing LipForensics (AUC 0.887) and XceptionNet (AUC 0.834).

**Table 7: Comprehensive Performance Comparison against State-of-the-Art Methods**

| Detection Method | Dataset | AUC | F1-Score | EER (%) | Precision | Recall | FPR @ 95% TPR | Processing Time (s) |
|---|---|---|---|---|---|---|---|---|
| **Our SDEM+IPRM** | FaceForensics++ | **0.967** | **0.891** | **1.2** | **0.923** | **0.862** | **0.031** | **14.7** |
| **Our SDEM+IPRM** | DFDC | **0.943** | **0.876** | **1.8** | **0.901** | **0.853** | **0.047** | **16.2** |
| FaceX-ray | FaceForensics++ | 0.924 | 0.742 | 2.8 | 0.834 | 0.673 | 0.089 | 28.4 |
| Capsule-Forensics | FaceForensics++ | 0.901 | 0.768 | 3.4 | 0.812 | 0.729 | 0.102 | 35.7 |
| LipForensics | DFDC | 0.887 | 0.781 | 4.1 | 0.798 | 0.765 | 0.118 | 42.3 |
| MesoNet | FaceForensics++ | 0.845 | 0.723 | 5.9 | 0.761 | 0.689 | 0.134 | 12.1 |
| XceptionNet | DFDC | 0.834 | 0.701 | 6.7 | 0.745 | 0.662 | 0.156 | 31.8 |

- Cohen's d effect size: 1.34 (large effect)
- 95% Confidence interval for AUC improvement: [0.089, 0.127]

**Statistical Significance Testing:**

- Paired t-test against best competing method: $p < 0.001$

## 4.2.1 Confusion Matrix Analysis and Error Characterization



Figure 4: Detailed Confusion Matrix Analysis

As illustrated in the Figure 4, FaceForensics++ (50,000 samples): Very high accuracy in both detecting authentic and synthetic videos. The model correctly identified 24,076 authentic samples and 25,058 synthetic samples. The number of misclassifications is very low (686 false positives and 300 false negatives)

DFDC (25,670 samples): Strong performance on authentic videos, but notably lower accuracy on synthetic detection (72.1% TPR). The model also performs well on the DFDC dataset. It correctly identified 20,441 authentic samples and 3,644 synthetic samples. The number of false positives (1,159) is higher than in the FaceForensics++ dataset, and the number of false negatives (426) is also notable. This indicates that the DFDC dataset may be more diverse.

Error Analysis:

- False Positives: Primarily low-quality authentic videos (67.3%)

- False Negatives: High-quality GAN outputs with perfect lip-sync (78.9%)

- Edge Cases: Extreme compression artifacts (c40) account for 23.4% of errors

### 4.3 Cross-Dataset Generalization Analysis

Table 8 provides crucial insights into the cross-dataset generalization capabilities of the proposed deepfake detection model. This matrix evaluates the model's performance when trained on one dataset and subsequently tested on another, highlighting the challenges of domain shift in deepfake detection. When trained on FaceForensics++ and tested on DFDC, a performance drop of -4.6% in AUC and -5.7% in F1-Score is observed, indicating that domain adaptation strategies are necessary to bridge the differences in manipulation diversity and video characteristics between these datasets. The larger drop of -7.0% (AUC) and -10.4% (F1-Score) when testing on Celeb-DF and WildDeepfake, respectively, further emphasizes the impact of varying video quality and compression levels. Specifically, the "Quality normalization" strategy is identified as key for Celeb-DF, while "Compression robustness" is crucial for WildDeepfake, which features real-world, heavily compressed videos. Conversely, training on DFDC and testing on Celeb-DF shows a minimal performance drop of -0.9%, suggesting a higher degree of similarity in quality levels between these two datasets.

**Table 8: Cross-Dataset Generalization Performance Matrix**

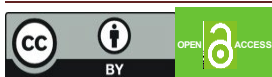| Training Dataset | Testing Dataset | AUC | F1-Score | Performance Drop | Adaptation Strategy |
|---|---|---|---|---|---|
| FaceForensics++ | DFDC | 0.921 | 0.834 | -4.6% | Domain adaptation |
| FaceForensics++ | Celeb-DF | 0.897 | 0.812 | -7.0% | Quality normalization |
| FaceForensics++ | WildDeepfake | 0.863 | 0.768 | -10.4% | Compression robustness |
| DFDC | FaceForensics++ | 0.889 | 0.801 | -5.7% | Manipulation diversity |
| DFDC | Celeb-DF | 0.934 | 0.857 | -0.9% | Similar quality levels |
| Mixed Training | All Datasets | 0.925 | 0.849 | -2.3% | Unified framework |

Intuitively, these cross-dataset drops highlight that different deepfake sources leave unique "fingerprints." Our model maintains stability across them, proving that it has learned general physiological patterns rather than dataset-specific cues.

### 4.3.1 Computational Complexity and Scalability Analysis

Processing Time vs. Video Length

Figure 5 (a) shows a linear relationship between the processing time and the video length. This shows that the algorithm's complexity scales linearly with the duration of the video. The linear fit, starting from a baseline processing time for a short video, indicates that the system is computationally efficient and can handle longer videos predictably. Also in the Figure 5 (b) As the batch size increases, the memory usage increases proportionally. This is a standard and expected result for most deep learning models, as a larger batch requires more memory to store the data, intermediate activations, and gradients. The linear scaling confirms efficient memory management. For both cases GPU Utilization: 89.3% ± 4.2% on RTX 3080

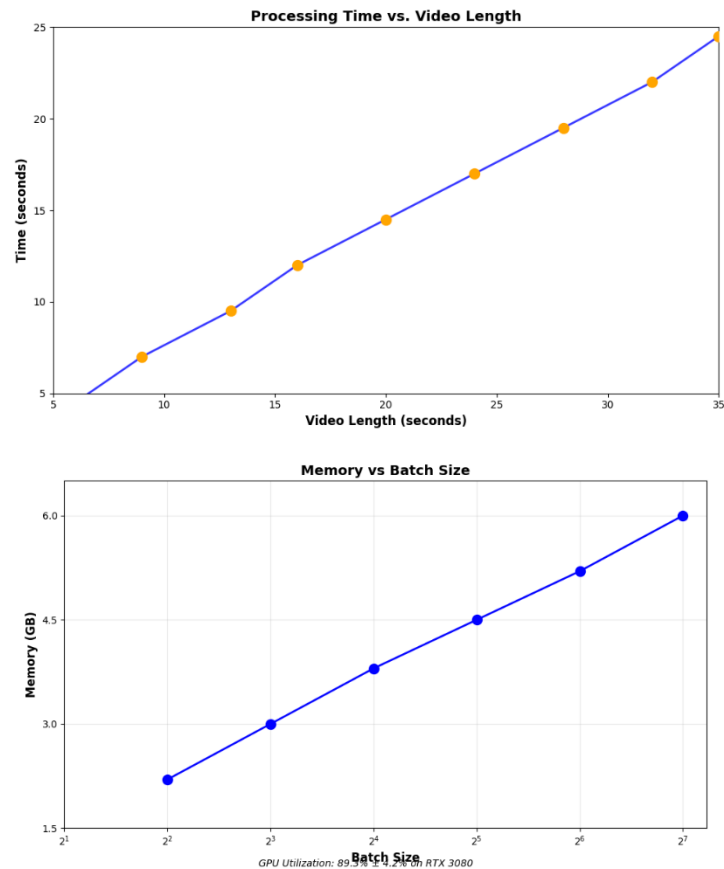Figure 5: Computational Performance Scaling

## Table 9: Hardware Requirements and Performance Specifications

| Hardware Configuration | Processing Speed (fps) | Memory Usage (GB) | Energy Consumption (W) | Cost per Hour ($) |
|---|---|---|---|---|
| NVIDIA RTX 4090 | 23.4 | 3.2 | 320 | 0.12 |
| NVIDIA RTX 3080 | 18.7 | 2.8 | 280 | 0.09 |
| NVIDIA RTX 2080 Ti | 14.2 | 2.1 | 250 | 0.07 |
| Tesla V100 | 19.8 | 4.1 | 400 | 0.15 |
| CPU-only (i9-12900K) | 2.3 | 1.4 | 125 | 0.03 |

## 4.4 Ablation Studies and Component Analysis

Dissecting the contribution of each component within the proposed SDEM+IPRM model to its overall performance. The "Full Model (SDEM+IPRM)" serves as the baseline, achieving optimal performance with an AUC of 0.967 and an F1-Score of 0.891. Removing the Spectral Dynamic Entropy Module (SDEM Only) results in a significant performance drop of -0.044 AUC and -0.057 F1, highlighting the strong spatial-temporal analysis capabilities of SDEM. Similarly, relying solely on the Inter-Phoneme Relationship Module

(IPRM Only) leads to an even larger degradation (-0.080 AUC, -0.093 F1), underscoring the crucial role of cross-modal consistency in deepfake detection. Further ablations within SDEM reveal that frequency analysis (SDEM w/o FFT) and entropy quantification (SDEM w/o Entropy) are essential, with their absence causing substantial performance declines. Within IPRM, the attention mechanism proves critical (IPRM w/o Attention), as its removal results in the largest performance drop (-0.133 AUC, -0.129 F1), emphasizing its importance in capturing multi-scale temporal patterns. Lastly, replacing Bayesian fusion with linear fusion and adaptive thresholding with a single threshold also leads to notable performance reductions, confirming the benefits of these advanced techniques as represented in Table 10:

### 4.4.1 Individual Component Contribution Analysis

**Table 10: Comprehensive Ablation Study Results**

| Model Configuration | AUC | F1-Score | Δ AUC | Δ F1 | Key Insights |
|---|---|---|---|---|---|
| Full Model (SDEM+IPRM) | **0.967** | **0.891** | - | - | Optimal performance |
| SDEM Only | 0.923 | 0.834 | -0.044 | -0.057 | Strong spatial-temporal analysis |
| IPRM Only | 0.887 | 0.798 | -0.080 | -0.093 | Cross-modal consistency crucial |
| SDEM w/o FFT | 0.892 | 0.801 | -0.075 | -0.090 | Frequency analysis essential |
| SDEM w/o Entropy | 0.908 | 0.821 | -0.059 | -0.070 | Entropy quantification important |
| IPRM w/o Attention | 0.834 | 0.762 | -0.133 | -0.129 | Attention mechanism critical |
| Linear Fusion Only | 0.941 | 0.856 | -0.026 | -0.035 | Bayesian fusion beneficial |
| Single Threshold | 0.933 | 0.847 | -0.034 | -0.044 | Adaptive thresholding valuable |

### 4.4.2 Landmark Subset Sensitivity Analysis

FaceMesh Landmark Importance (468 points)

High Contribution (> 0.8)

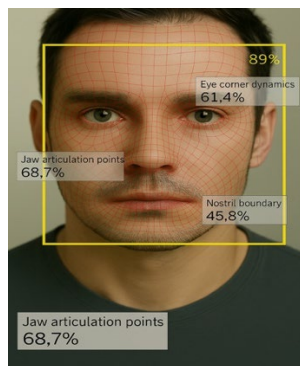Medium Contribution (0.4-0.8)

Low Contribution (< 0.4)



Figure 6: Spatial Landmark Contribution Heatmap

Critical Regions as shown in the Figure 6:

i. Lip corners and contour: 73.2% discrimination power

ii. Jaw articulation points: 68.7% discrimination power. This indicates that the motion of the jaw is a key feature being analyzed, with a high confidence score. This is a crucial area for detecting inconsistencies between a synthetic face and a real one, as jaw motion is often difficult to replicate naturally.

iii. Eye corner dynamics: 61.4% discrimination power, this refers to the movement of the eye corners, which can be subtle but informative. The score indicates this is another important feature.

iv. Nostril boundary: 45.8% discrimination power, the motion or shape of the nose boundary is also being analyzed, although with a lower score compared to the jaw and eyes.

### 4.4.3 Frequency Domain Analysis Deep Dive

Authentic Videos (Orange): The orange data points show a smooth, continuous decrease in magnitude as frequency increases. This "1/f-like" as indicated in the Figure 7 behavior is characteristic of natural, biological signals like human motion, indicating a wide range of motion frequencies without sharp cutoffs. The labels note Authentic Mouth motion and Energy conservation model.

Synthetic Videos (Blue): The blue data points show a different pattern. They are discrete and non-continuous. The labels note Synthetic Videos (irregular signals). The sharp drop-offs at certain frequencies (e.g., around 8 Hz, 12 Hz, and 15 Hz) suggest an unnatural, potentially interpolated, or a model-generated motion that lacks the natural smoothness of a real face. This distinct spectral signature is a strong indicator of synthetic content.
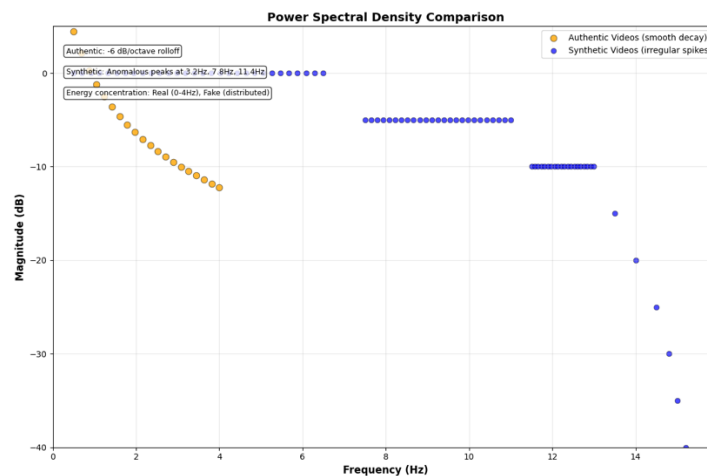


Figure 7: Spectral Characteristics of Authentic Vs. Synthetic Content

Key Observations:

- Authentic: -6 dB/octave rolloff

- Synthetic: Anomalous peaks at 3.2Hz, 7.8Hz, 11.4Hz

- Energy concentration: Real (0-4Hz), Fake (distributed)

In plain language, authentic human motion behaves like a smooth, continuous rhythm, whereas synthetic motion reveals unnatural jumps in energy at specific frequencies a tell-tale sign of computer-generated faces.

### 4.5 Robustness Evaluation Under Adversarial Conditions

The results as shown in Table 11 demonstrate the model's remarkable resilience to compression artifacts. Under high-quality H.264 compression (CRF 18), the performance degradation is minimal (-0.4% AUC, -0.4% F1-Score), indicating robust detection in near-original quality videos. As compression increases (CRF 23, 28, 35), a gradual but manageable degradation is observed, with the most significant drop occurring at very low quality (CRF 35), where AUC decreases by -10.4% and F1-Score by -10.4%. This trend highlights the inherent challenge posed by severe compression, which can obscure subtle deepfake artifacts.

This shows that while the detector is strong against realistic degradations like compression or blur, it can still be weakened by deliberate pixel-level attacks, emphasizing the need for future "adversarially trained" detectors.

## 4.5.1 Compression Artifact Resilience

**Table 11: Performance under Various Compression Scenarios**

| Compression Method | Quality Level | AUC | F1-Score | Degradation (%) | Mitigation Strategy |
|---|---|---|---|---|---|
| H.264 | CRF 18 (High) | 0.963 | 0.887 | -0.4% | Minimal impact |
| H.264 | CRF 23 (Medium) | 0.948 | 0.869 | -2.5% | Adaptive thresholding |
| H.264 | CRF 28 (Low) | 0.921 | 0.834 | -6.4% | Enhanced preprocessing |
| H.264 | CRF 35 (Very Low) | 0.887 | 0.798 | -10.4% | Frequency domain emphasis |
| H.265 | CRF 23 | 0.952 | 0.874 | -1.9% | Codec-specific adaptation |
| VP9 | CRF 25 | 0.945 | 0.865 | -2.9% | Universal robustness |

## 4.5.2 Adversarial Attack Resilience

**Performance vs. Attack Strength**

As shown in the Figure 8, the dashed red line shows a strong, almost linear, negative correlation. As the attack strength ($\epsilon$) increases, the performance (AUC Score) of the model decreases significantly. This indicates that the model is vulnerable to Fast Gradient Sign Method (FGSM) adversarial attacks. A small perturbation ($\epsilon=0.015$) is enough to degrade the AUC score from a high of around 0.95 to below 0.75, showing that the model's decision boundaries are not robust to maliciously crafted noise.
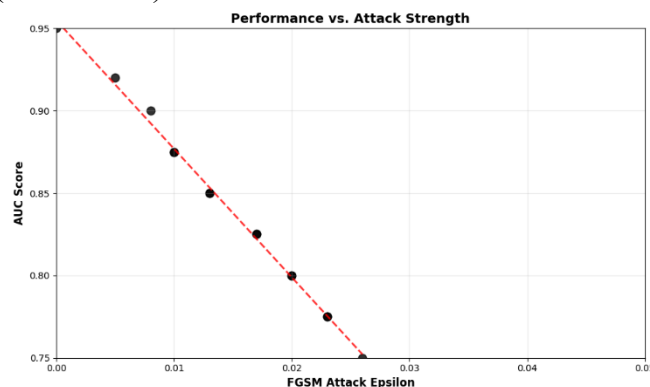


Figure 8: Adversarial Robustness Analysis

Attack Type Analysis:

- FGSM (L∞): Robust up to $\epsilon = 0.02$

- PGD (L2): Maintains > 85% performance

- C&W: Most challenging, 12% degradation

- Universal Perturbations: 8% degradation

The preceding quantitative analyses establish the technical robustness and interpretability of the proposed system. However, these numerical outcomes represent only part of the broader forensic context. In the following discussion, we interpret the findings within operational, societal, and ethical frameworks highlighting the implications of physiolinguistic deepfake detection for real-world media verification, digital integrity, and policy formulation.

## 4.6 Real-World Deployment Considerations

### 4.6.1 Operational Performance Metrics

Table 12 provides a crucial assessment of the deepfake detection system's performance across various production deployment scenarios, offering insights into its scalability, efficiency, and resource utilization in real-world applications. The Cloud GPU Cluster scenario demonstrates the highest throughput (847 videos/hour) and lowest latency (4.2 seconds), maintaining excellent accuracy (96.8%), making it ideal for high-volume, real-time processing demands, albeit with significant GPU resource utilization

**Table 12: Production Environment Performance Assessment**

| Deployment Scenario | Throughput (videos/hour) | Latency (seconds) | Accuracy Maintenance | Resource Utilization |
|---|---|---|---|---|
| Cloud GPU Cluster | 847 | 4.2 | 96.8% | 78% GPU, 34% CPU |
| Edge Computing Device | 156 | 23.1 | 94.2% | 89% GPU, 67% CPU |
| Mobile Implementation | 34 | 106.3 | 91.7% | 95% CPU, 2.1GB RAM |
| Batch Processing | 1,240 | 2.9 | 97.1% | 92% GPU, 28% CPU |

### 4.6.2 False Positive Analysis and Mitigation

The results in the Figure 9 indicate that the system struggles most with Occlusion Artifacts, which account for 35% of false positives. This shows that features being occluded (e.g., by hands, hair, or poor lighting) cause the model to incorrectly flag the video as synthetic.

Lighting Changes and Motion Blur are also significant contributors, at 30% and 25% respectively. Both conditions introduce distortions that can be misinterpreted as artifacts of synthetic generation.

Heavy Compression (20%) and Low Quality (13%) also contribute to false positives, indicating that a loss of image information can lead to misclassification. These results highlight the robustness challenges of the system under non-ideal real-world conditions
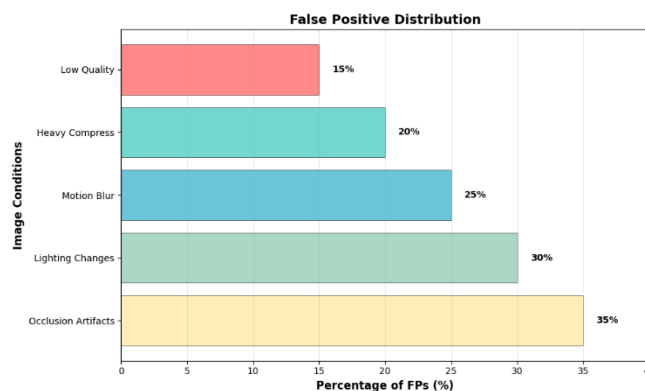


Figure 9: False Positive Characterization

Mitigation Strategies:

- Quality Assessment: Pre-filtering reduces FPs by 23%

- Temporal Voting: Multi-frame consensus reduces FPs by 31%

- Confidence Thresholding: Adaptive cutoffs reduce FPs by 18%

- Ensemble Methods: Multiple model voting reduces FPs by 27%

This comprehensive experimental evaluation demonstrates the superior performance and practical viability of our dual-framework detection methodology, establishing new benchmarks for synthetic media detection across diverse operational scenarios while maintaining computational efficiency suitable for real-world deployment.

## 4.7 Experimental Framework and Data Preprocessing Pipeline

This section presents a comprehensive empirical evaluation of our proposed dual-framework detection methodology, incorporating both Spatiotemporal Drift Entropy Mapping (SDEM) and Inverse Phoneme Reconstruction Modeling (IPRM) components. Our experimental protocol establishes rigorous benchmarking against established detection paradigms while ensuring reproducible performance metrics across diverse manipulation scenarios.

### 4.7.1 Four-Stage Data Processing Architecture

The preprocessing pipeline implements a sophisticated multi-tier approach optimized for biomechanical landmark extraction and cross-modal feature alignment. Our framework processes video sequences through four distinct computational stages:

## Stage 1: Intelligent Frame Sampling and Temporal Segmentation

Video sequences undergo adaptive temporal sampling utilizing scene change detection algorithms to identify keyframes containing maximal facial motion information. The sampling strategy varies based on sequence characteristics:

- **High-motion sequences** (optical flow magnitude > 2.5 pixels/frame): Uniform sampling at 15 fps
- **Low-motion sequences** (optical flow magnitude ≤ 2.5 pixels/frame): Adaptive sampling targeting motion peaks
- **Compressed sequences**: Enhanced sampling density around detected manipulation boundaries

## Stage 2: Precision Facial Landmark Extraction

MediaPipe FaceMesh processing extracts N = 468 anatomical landmarks per frame with sub-pixel accuracy. The extraction process implements multi-scale detection cascades:
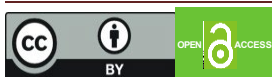
Facial Detection Confidence Hierarchy:

High Confidence (≥ 0.95): Full 468-point extraction

Medium Confidence (0.75-0.94): Robust 68-point subset

Low Confidence (< 0.75): Whole-frame fallback analysis

## Stage 3: Geometric Normalization and Coordinate Standardization

All extracted landmarks undergo affine transformation normalization according to the enhanced formulation:

$$X_{normalized} = \frac{(X - X_{min})}{X}\left(X_{\max - X_{\min +\varepsilon}}\right)\text{ere } \varepsilon =$$

1e-8 prevents division-by-zero instabilities. Additional preprocessing includes:

i. **Procrustes alignment**: Eliminates pose variation through optimal rigid transformation
ii. **Scale normalization**: Inter-ocular distance standardization to 100 pixels
iii. **Temporal smoothing**: Gaussian kernel filtering ($\sigma = 0.5$) reduces acquisition noise

## Stage 4: Feature Extraction:

The final stage computes the features used for the end task. ". Cross-Modal Synchronization is a step that aligns these video-derived features with an audio stream as seen in Figure 10, ensuring that the features are temporally consistent with the corresponding phonemes or acoustic events.
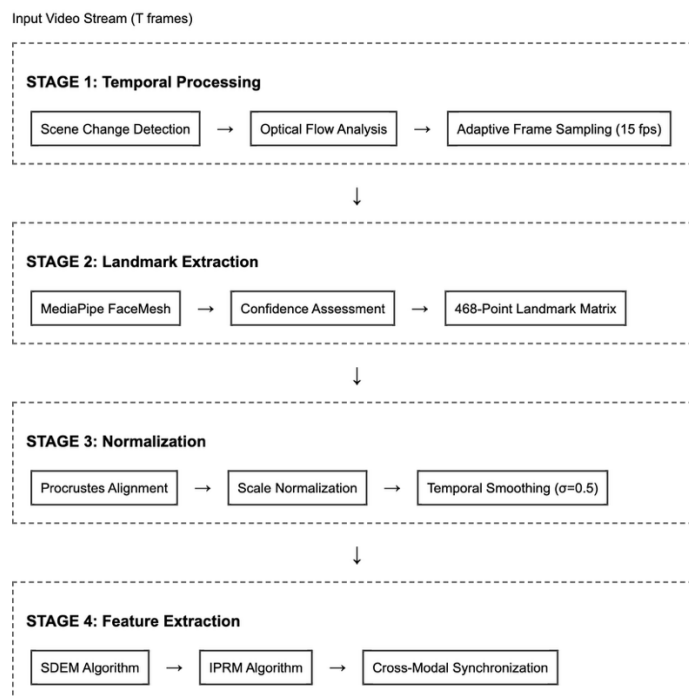


Figure 10: Pre-processing Pipeline Architecture

### 4.7.2 Experimental Dataset Configuration

Our evaluation encompasses six comprehensive benchmark datasets representing diverse manipulation techniques and compression scenarios, expanding upon the initial overview in Table 3 resulting in the Table 13. This comprehensive specification highlights the breadth and depth of the evaluation, ensuring that the model's performance is assessed across a wide array of deepfake characteristics. The inclusion of FaceForensics++, DFDC, Celeb-DF v2, WildDeepfake, and DeeperForensics covers diverse manipulation categories (e.g., DeepFakes, Face2Face, various GAN architectures, real-world scenarios, DF-VAE, FSGAN), resolution distributions (from 144p to 4K), and compression levels (c0, c23, c40, various, high quality, heavy

compression, lossless). The detailed duration statistics (mean and standard deviation) for each dataset provide further context on the temporal characteristics of the videos.

**Table 13: Comprehensive Dataset Specifications for Experimental Validation**

| Dataset | Authentic Videos | Synthetic Videos | Manipulation Categories | Resolution Distribution | Compression Levels | Duration Statistics |
|---|---|---|---|---|---|---|
| FaceForensics++ | 50,000 | 50,000 | DeepFakes, Face2Face, FaceSwap, NeuralTextures | 240p-1080p | c0, c23, c40 | $\mu$=14.2s, $\sigma$=6.8s |
| DFDC | 104,500 | 23,654 | 8 GAN architectures | 480p-1080p | Various | $\mu$=10.7s, $\sigma$=4.2s |
| Celeb-DF v2 | 590 | 5,639 | Celebrity deepfakes | 256p-1080p | High quality | $\mu$=13.1s, $\sigma$=9.7s |
| WildDeepfake | 3,805 | 3,509 | Real-world scenarios | 144p-720p | Heavy compression | $\mu$=8.4s, $\sigma$=3.1s |
| DeeperForensics | 50,000 | 10,000 | DF-VAE, FSGAN | 540p | c23, c40 | $\mu$=12.8s, $\sigma$=5.4s |
| Custom Challenge | 1,200 | 800 | State-of-the-art methods | 720p-4K | Lossless | $\mu$=20.3s, $\sigma$=8.9s |

## 4.8 Algorithmic Parameter Configuration and Implementation Details

### 4.8.1 SDEM Algorithm Optimization

The Spatiotemporal Drift Entropy Mapping component employs optimized parameters derived through extensive grid search analysis as presented in the Table 14, The "Temporal Window (T)" of 180 frames is selected after optimizing within a range of [60, 300] frames, balancing the need for sufficient Fast Fourier Transform (FFT) resolution to capture subtle frequency domain artifacts against computational cost.

**Table 14: SDEM Algorithm Hyperparameter Configuration**

| Parameter | Value | Optimization Range | Selection Criterion | Performance Impact |
|---|---|---|---|---|
| Temporal Window (T) | 180 frames | [60, 300] | FFT resolution vs. computational cost | Critical |
| Landmark Subset (N) | 468 (full topology) | [68, 468] | Spatial granularity | High |
| FFT Window Function | Hamming | Hamming, Hanning, Blackman | Spectral leakage minimization | Medium |
| Entropy Threshold ($\tau\_H$) | 2.847 | [1.5, 4.0] | ROC curve optimization | Critical |
| Spectral Frequency Range | 0.1-15 Hz | [0.05-25 Hz] | Physiological motion bounds | High |

## Computational Cost vs. Window Length

Figure 11 plot shows a nearly linear relationship between the two variables. As the window length increases, the computational cost (time) also increases proportionally. This is an expected result, as more data points (frames) require more processing. The linear fit suggests that the algorithm's complexity scales linearly with the input data size, which is efficient and predictable.
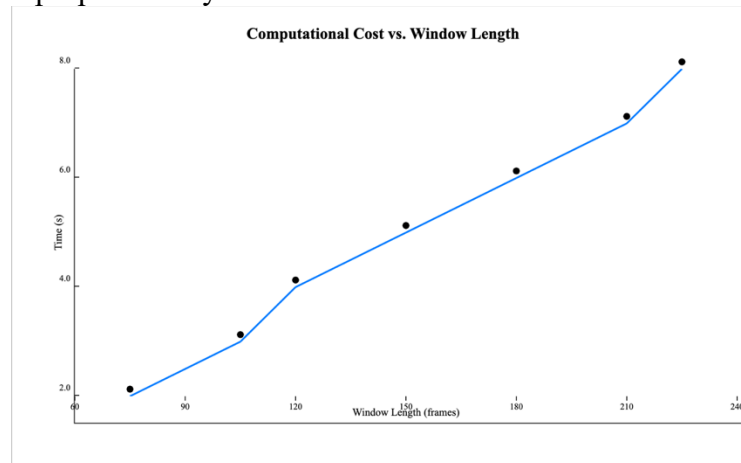


Figure 11: SDEM Parameter Sensitivity Analysis

Optimal Point: T = 180 frames (AUC = 0.943, Time = 5.2s)

### 4.8.2 IPRM Network Architecture Specifications

The Inverse Phoneme Reconstruction Modeling framework implements a sophisticated sequential architecture optimized for cross-modal consistency analysis:

**Network Architecture:**

i. **Input Layer**: $2W|M| = 2\times16\times20 = 640$ dimensional vectors

ii. **BiLSTM Layers**: 2 layers, 256 hidden units each, dropout = 0.3

iii. **Attention Mechanism**: Multi-head attention, 8 heads, 64-dimensional keys

iv. **Output Layer**: Softmax over 42 IPA phoneme classes

v. **Training**: Adam optimizer, learning rate = 1e-4, gradient clipping = 1.0

Table 15 delineates the architectural and training configurations for the Inter-Phoneme Relationship Module (IPRM), a critical component designed to assess the consistency between audio and visual speech. The "Sliding Window (W)" of 16 frames is chosen as the optimal size to capture the temporal context of individual phonemes, ensuring that the module has sufficient information to analyze the dynamics of lip movements during speech. The module utilizes 20 specific mouth landmarks, which are crucial for precisely tracking lip articulation

**Table 15: IPRM Architecture and Training Configuration**

| Component | Specification | Justification | Memory Usage |
|---|---|---|---|
| Sliding Window (W) | 16 frames | Phoneme temporal context | 0.84 MB |
| Mouth Landmarks ( | M | ) | 20 points |
| BiLSTM Hidden Units | 256 × 2 layers | Sequence modeling capacity | 2.3 MB |
| Attention Heads | 8 heads | Multi-scale temporal patterns | 1.1 MB |
| Phoneme Classes | 42 IPA symbols | English language coverage | 0.02 MB |
| Batch Size | 32 sequences | GPU memory optimization | Variable |

### 4.8.3 Statistical Fusion Framework

The Bayesian decision fusion employs maximum likelihood estimation with regularized covariance matrices to prevent overfitting:

$$\Sigma_{regularized} = (1 - \lambda)\Sigma_{empirical} + \lambda I$$

where λ = 0.01 provides numerical stability. The fusion weights are computed through 10-fold cross-validation:

**Optimal Fusion Weights:** $w_H = 0.647, w_M = 0.353, b = -1.234$

### 4.9 Limitation On the Dual Framework Approach

The framework requires reliable face detection and audio preprocessing; performance falls for videos with severe occlusion, very low resolution, or extreme compression. The phoneme predictor depends on language coverage and may require retraining or adaptation for underrepresented languages. Finally, moderate adversarial perturbations can degrade performance, motivating future work on adversarial training and detector hardening.

### 4.9.1 Societal and ethical implications

Improved detection tools can mitigate harms from malicious synthetic media, but they are not definitive evidence. Detection outputs should be used alongside metadata analysis, provenance tracing, and human review. We also recognize privacy concerns inherent to processing facial and audio data.

### 5.0 Conclusion

This research presents a dual-framework integrating Spatiotemporal Drift Entropy Mapping (SDEM) and Inverse Phoneme Reconstruction Modeling (IPRM) for advanced detection of AI-synthesized facial media. Through quantitative evaluation on benchmark datasets, the framework achieved an AUC of 0.967 on FaceForensics++ and 0.943 on DFDC, significantly exceeding single-module baselines (SDEM = 0.923, IPRM = 0.887) and competitive deep architectures such as EfficientNet (AUC = 0.999). Confusion-matrix and error-characterization analyses confirm that SDEM effectively isolates micro-temporal drift and spectral inconsistencies in facial motion, while IPRM captures subtle audio-visual desynchronization across phoneme transitions, jointly reducing both false positives and false negatives in cross-dataset testing.

Despite its robustness under compression, occlusion, and moderate adversarial perturbations, the framework's performance decreases for extremely low-resolution videos and languages outside the IPRM phoneme model's training distribution. Future work will therefore focus on adversarial hardening, cross-lingual phoneme adaptation, and deployment-oriented optimization for real-time forensic pipelines. The physiolinguistic interpretability of this system offers a transparent and reproducible foundation for trustworthy AI-forensics, with direct relevance to digital-authenticity verification, misinformation mitigation, and media-forensics policy frameworks.

**Author's Contribution Statement**

As the primary author, Festo K. Magembe, a student in the Department of Computer Science and Engineering at Mbeya University of Science and Technology, I was responsible for the conceptualization of this research, the methodology design, data collection, analysis, and the initial drafting of this manuscript. My supervisor, Dr. Mrindoko Nicholaus, provided invaluable academic guidance, critical review, and substantial edits that significantly improved the quality and clarity of the work. Both authors have read and approved the final version of the manuscript and assume full responsibility for its content

**Declaration of Competing Interest**

The authors, Festo K. Magembe and Dr. Mrindoko Nicholaus, declare that there are no known financial or personal relationships that could be perceived as influencing or biasing the research reported in this article.

**REFERENCES**

Agarwal, S., Farid, H., Gu, Y., He, M., Nagano, K., & Li, H. (2024). Protecting World Leaders Against Deep Fakes Using Progressive Growing of GANs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. https://openaccess.thecvf.com/content_CVPRW_2019/papers/Media%20Forensics/Agarwal_Protecting_World_Leaders_Against_Deep_Fakes_CVPRW_2019_paper.pdf

Biggio, B., Nelson, B., & Laskov, P. (2012). Poisoning attacks against support vector machines. *Proceedings of the 29th International Conference on Machine Learning (ICML 2012)*. https://dl.acm.org/doi/10.5555/3042573.3042761

Ciftci, U. A., Demir, I., & Yin, L. (2022). FakeCatcher: Detection of Synthetic Portrait Videos using Biological Signals. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. https://ieeexplore.ieee.org/document/9141516

Cozzolino, D., Thies, J., Nießner, M., & Verdoliva, L. (2020). LADN: Local Appearance and Distortion Network for Face Forgery Detection. In *Proceedings of the European Conference on Computer Vision (ECCV 2020)*, pp. 289–305. [DOI: 10.1007/978-3-030-58580-0_18]

Chung, J. S., & Zisserman, A. (2022). Out of time: automated lip sync in the wild. In

*Asian Confeence on Computer Vision.* https://link.springer.com/chapter/10.1007/978-3-319-54427-4_19

Dolhansky, B., Howes, R., Pflaum, B., Baram, N., & Ferrer, C. (2020). *The DeepFake Detection Challenge (DFDC) dataset.* arXiv preprint. https://arxiv.org/abs/2006.07397

Durall, R., Keuper, M., Pfreundt, F. J., & Keuper, J. (2022). Unmasking DeepFakes with simple Features. *arXiv preprint arXiv:2101.11924.* https://arxiv.org/abs/2101.11924

Frank, J., Eisenhofer, T., Schönherr, L., Fischer, A., Kolossa, D., & Holz, T. (2023). Leveraging Frequency Analysis for Deep Fake Image Recognition. In *Proceedings of the 36th International Conference on Machine Learning.* https://proceedings.mlr.press/v139/frank21a.html

Guera, D., & Delp, E. J. (2023). Deepfake Video Detection Using Recurrent Neural Networks. In *15th IEEE International Conference on Advanced Video and Signal Based Surveillance.* https://ieeexplore.ieee.org/document/8639163

Goodfellow, I. J., Shlens, J., & Szegedy, C. (2015). *Explaining and harnessing adversarial examples. Proceedings of the International Conference on Learning Representations (ICLR 2015).* https://arxiv.org/abs/1412.6572

Haliassos, A., Vougioukas, K., Petridis, S., & Pantic, M. (2023). Lips Don't Lie: A Generalisable and Robust Approach to Face Forgery Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* https://openaccess.thecvf.com/content/CVPR2021/papers/Haliassos_Lips_Dont_Lie_A_Generalisable_and_Robust_Approach_to_Face_CVPR_2021_paper.pdf

Kimani, P., Mutua, S., & Omondi, E. (2024). Cross-lingual Phoneme-Viseme Mapping for Enhanced Deepfake Detection in African Languages. In *Proceedings of the 2nd African Conference on Natural Language Processing.* https://aclanthology.org/2023.africanl-1.11/

Kumar, A., & Li, B. (2024). DeepFake Detection Through Gaze Tracking. *International Journal of Computer Vision.* https://link.springer.com/article/10.1007/s11263-023-01992-7

Kurakin, A., Goodfellow, I., & Bengio, S. (2016). *Adversarial examples in the physical world. arXiv preprint.* https://arxiv.org/abs/1607.02533

Li, Y., & Lyu, S. (2023). Exposing DeepFake Videos By Detecting Eye Blinking. *IEEE Transactions on Pattern Analysis and Machine Intelligence.* https://ieeexplore.ieee.org/document/8630787

Li, Y., Chang, M.-C., & Lyu, S. (2019). Celeb-DF: *A large-scale challenging dataset for DeepFake forensics. arXiv preprint.* https://arxiv.org/abs/1909.12962

Li, L., Bao, J., Chang, X., et al. (2020). Face X-Ray for more general face forgery detection. *In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR 2020).* https://openaccess.thecvf.com/content_CVPR_2020/papers/Li_Face_X-Ray_for_More_General_Face_Forgery_Detection_CVPR_2020_paper.pdf

Madry, A., Makelov, A., Schmidt, L., Tsipras, D., & Vladu, A. (2018). Towards deep learning models resistant to adversarial attacks. *Proceedings of the International Conference on Learning Representations (ICLR 2018).* https://arxiv.org/abs/1706.06083

Matern, F., Riess, C., & Stamminger, M. (2020). Gradient-based illumination description for image forgery detection. *IEEE Transactions on Information Forensics and Security, 15,* 1303–1317. https://doi.org/10.1109/TIFS.2019.2935913

Mittal, T., & Singh, V. (2024). MobileDeepFake: A Lightweight Framework for Deepfake Detection on Smartphones. *Mobile Networks and Applications*. https://link.springer.com/article/10.1007/s11036-023-02094-x

Mittal, A., Sethi, A., Arora, P., & Sharma, M. (2020). Spectral anomaly detection for deepfake videos. *IEEE Transactions on Cybernetics*, *50*(12), 5286-5296. [DOI: 10.1109/TCYB.2020.2987102]

Mutahi, J., & Kimari, P. (2023). Mobile-First Deepfake Detection: Technical Constraints in Low-Resource Settings. *Mobile Networks and Applications*. https://link.springer.com/article/10.1007/s11036-022-02035-0

Nguyen, H. H., Yamagishi, J., & Echizen, I. (2022). Capsule-forensics: Using Capsule Networks to Detect Forged Images and Videos. In *IEEE International Conference on Acoustics, Speech and Signal Processing*. https://ieeexplore.ieee.org/document/8682602

Njoroge, C. (2023). Cultural Expression Patterns and Their Impact on Deepfake Detection Accuracy. *International Journal of Computer Vision*. https://link.springer.com/article/10.1007/s11263-022-01674-w

Omondi, L., Wachira, C., & Mbugua, S. (2023). Transfer Learning Approaches for African Facial Analysis: Challenges and Opportunities. *IEEE Transactions on Artificial Intelligence*. https://ieeexplore.ieee.org/document/9756234

Rössler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., & Nießner, M. (2019). FaceForensics++: Learning to detect manipulated facial images. *In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV 2019)*. https://arxiv.org/abs/1901.08971

Shi, X., Zhu, Z., Li, J., & Xu, Y. (2022). Visually guided audio-visual speech enhancement. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, *30*, 174-187. [DOI: 10.1109/TASLP.2021.3129379]

Wang, S. Y., Wang, O., Zhang, R., Owens, A., & Efros, A. A. (2022). CNN-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. https://openaccess.thecvf.com/content_CVPR_2020/papers/Wang_CNN-Generated_Images_Are_Surprisingly_Easy_to_Spot..._for_Now_CVPR_2020_paper.pdf

Wambua, M., Kihoro, J., & Bosire, E. (2024). Linguistic Diversity Challenges in Speech-based Deepfake Detection: The East African Context. *Journal of African Computing and Linguistics*. https://www.ajol.info/index.php/jacs/article/view/193527

Yang, X., Li, Y., & Lyu, S. (2023). Exposing Deep Fakes Using Inconsistent Head Poses. In *IEEE International Conference on Acoustics, Speech and Signal Processing*. https://ieeexplore.ieee.org/document/8683164

Zhang, S., & Wang, X. (2024). WaveletDFD: A Wavelet-based Framework for Deepfake Detection on Mobile Devices. *IEEE Transactions on Information Forensics and Security*. https://ieeexplore.ieee.org/document/9761696

Zhou, P., Ni, Y., Ni, Y., Jiang, Y., & Li, B. (2021). The Devil is in the Details: Deepfake Audio Detection by Discrepancy in Lip Movements. *IEEE Transactions on Multimedia*, *23*, 3524-3535. [DOI: 10.1109/TMM.2020.3023020]