

Adaptive RGB-Based Face Spoof Detection Using Multi-Expert Feature Decomposition and Gated Fusion

Vikas Kumar

Department of Computer Engineering
Army Institute of Technology
Pune, Maharashtra, India 411015
Vikaskumar_240252@aitpune.edu.in

Abstract

Face spoofing attacks—including print, replay, and digitally manipulated facial presentations—pose a significant threat to modern biometric authentication systems. Conventional spoof detection approaches often rely on single-stream convolutional models or specialized hardware sensors such as depth or infrared cameras, limiting their robustness and deployability in real-world consumer devices. To address these challenges, this paper proposes an adaptive multi-expert spoof detection framework that operates exclusively on a single RGB input while achieving enhanced generalization against both known and unseen attacks. The architecture decomposes spoof detection into three complementary expert branches: an RGB appearance expert, a depth-aware structural expert, and a frequency-domain expert. To effectively integrate these heterogeneous cues, a reliability-aware gated fusion mechanism is introduced to dynamically weight expert contributions. Furthermore, the framework incorporates an adaptive feedback mechanism that continuously refines the model by leveraging high-confidence samples during deployment. Extensive experiments demonstrate that the proposed approach achieves superior performance across diverse attack scenarios compared to state-of-the-art methods.

Index Terms

Face Spoof Detection, Presentation Attack Detection, RGB-Based Biometrics, Multi-Expert Learning, Gated Fusion, Frequency Domain Analysis, Adaptive Learning.

I. INTRODUCTION

Face recognition systems are widely deployed in security-critical applications such as mobile authentication, access control, and financial services. Despite their success, these systems remain highly vulnerable to Presentation Attacks (PAs), also known as face spoofing. Common attacks include printed photographs, replayed videos on high-definition screens, and 3D masks. Such attacks can be performed with minimal resources and pose a serious threat to the reliability of biometric authentication systems.

To address this challenge, face spoof detection (or Presentation Attack Detection - PAD) has become an essential component of modern biometric pipelines. Early approaches relied on handcrafted texture descriptors (e.g., LBP, HOG), while recent methods predominantly employ deep learning models trained on RGB images. However, many existing solutions adopt single-stream convolutional architectures, implicitly assuming that a single representation is sufficient to capture all spoofing artifacts.

In practice, spoofing attacks manifest across multiple perceptual dimensions:

- **Appearance:** Color distortion and texture loss.
- **Geometry:** Lack of 3D volumetric structure in flat media.
- **Spectral:** Moiré patterns and high-frequency noise in digital displays.

While several works have explored using depth or infrared sensors to detect these cues, such hardware is often unavailable on consumer-grade devices. Furthermore, spoofing strategies evolve rapidly, leading to “domain shift” where models trained on fixed datasets fail against unseen attack types.

Motivated by these limitations, this paper proposes an **Adaptive RGB-Based Face Spoof Detection Framework**. We explicitly model spoof detection as a multi-expert perceptual reasoning problem. A single RGB input is processed by specialized expert branches focusing on appearance, structural depth, and frequency analysis. Crucially, we introduce a **Reliability-Aware Gated Fusion** mechanism that dynamically estimates the confidence of each expert to weigh their contributions. Finally, to handle evolving threats, we implement an **Adaptive Feedback Mechanism** (Fig. 1) that updates the model using a replay buffer of high-confidence samples collected during inference.

II. RELATED WORK

Existing PAD approaches can be broadly categorized into appearance-based, depth-based, and frequency-domain methods.

A. Appearance-Based Spoof Detection

Early methods relied on handcrafted features like Local Binary Patterns (LBP) and Histogram of Oriented Gradients (HOG) to capture texture anomalies. While computationally efficient, these methods struggle with varying illumination and high-quality replay attacks. Deep Learning approaches, particularly Convolutional Neural Networks (CNNs), have since surpassed handcrafted methods by learning discriminative features directly from data. However, standard CNNs often overfit to dataset-specific background cues rather than the spoofing artifacts themselves.

B. Depth and Structural Cue-Based Methods

To overcome RGB limitations, researchers have incorporated depth sensors or structured light. These methods effectively detect flat attacks (e.g., photos) by analyzing surface geometry. In the absence of depth sensors, recent works estimate “pseudo-depth” maps from RGB images using deep supervision. Our work adopts this RGB-to-Depth strategy but treats it as a distinct “expert” branch rather than just an auxiliary loss.

C. Frequency-Domain Approaches

Frequency analysis is effective for detecting artifacts introduced by resizing, printing, or screen refresh rates (Moiré patterns). Fourier Transform (FT) based methods analyze high-frequency components that are invisible in the spatial domain. Recent hybrid models combine spatial and frequency features, but often use static concatenation, ignoring that frequency cues may be unreliable in low-resolution images.

D. Multi-Modal and Adaptive Fusion

Fusion strategies typically involve simple averaging or concatenation. Few works employ attention mechanisms to weigh modalities dynamically. Moreover, most existing PAD systems are static; they do not adapt after deployment. Our work distinguishes itself by integrating dynamic gated fusion with an online adaptive feedback loop.

III. SYSTEM ARCHITECTURE

The proposed framework is designed as a modular, multi-stream architecture that processes a single RGB image to distinguish between *bona fide* (live) and *presentation attack* (spoof) faces. As illustrated in Fig. 1, the pipeline consists of four primary modules: (A) Heterogeneous Expert Branches for multi-domain feature extraction, (B) A Reliability-Aware Gated Fusion module for dynamic feature integration, (C) A Decision & Deployment module for final classification, and (D) An Adaptive Feedback Mechanism for continuous model refinement.

A. Heterogeneous Expert Branches

To capture the diverse manifestations of spoofing attacks, the input RGB image I_{rgb} is processed in parallel by three specialized expert branches. Each branch focuses on a distinct perceptual domain:

1) *RGB Appearance Branch*: This branch targets spatial domain artifacts such as surface texture anomalies, color distortion, and screen moiré patterns. It utilizes a lightweight Convolutional Neural Network (CNN) backbone to extract a high-level appearance embedding, denoted as F_{rgb} . This expert is particularly effective at identifying low-quality replay attacks and printed photo artifacts visible in the visual spectrum.

2) *Depth-Aware Structural Branch*: Since standard RGB cameras lack depth sensing capabilities, this branch employs a depth-estimation encoder-decoder network to predict a pseudo-depth map from the single RGB input. The network is trained with auxiliary supervision to distinguish between the volumetric surface of a live face and the planar surface of a presentation attack (e.g., a photo or screen). The latent representation from the depth encoder serves as the structural feature vector, F_{depth} .

3) *Frequency Domain Branch*: To detect high-frequency anomalies introduced by digital resizing or printing processes (e.g., blurring or grid artifacts), the input image is transformed into the spectral domain via Fast Fourier Transform (FFT). The magnitude spectrum is processed by a shallow CNN to extract spectral features, yielding the frequency domain embedding F_{freq} .

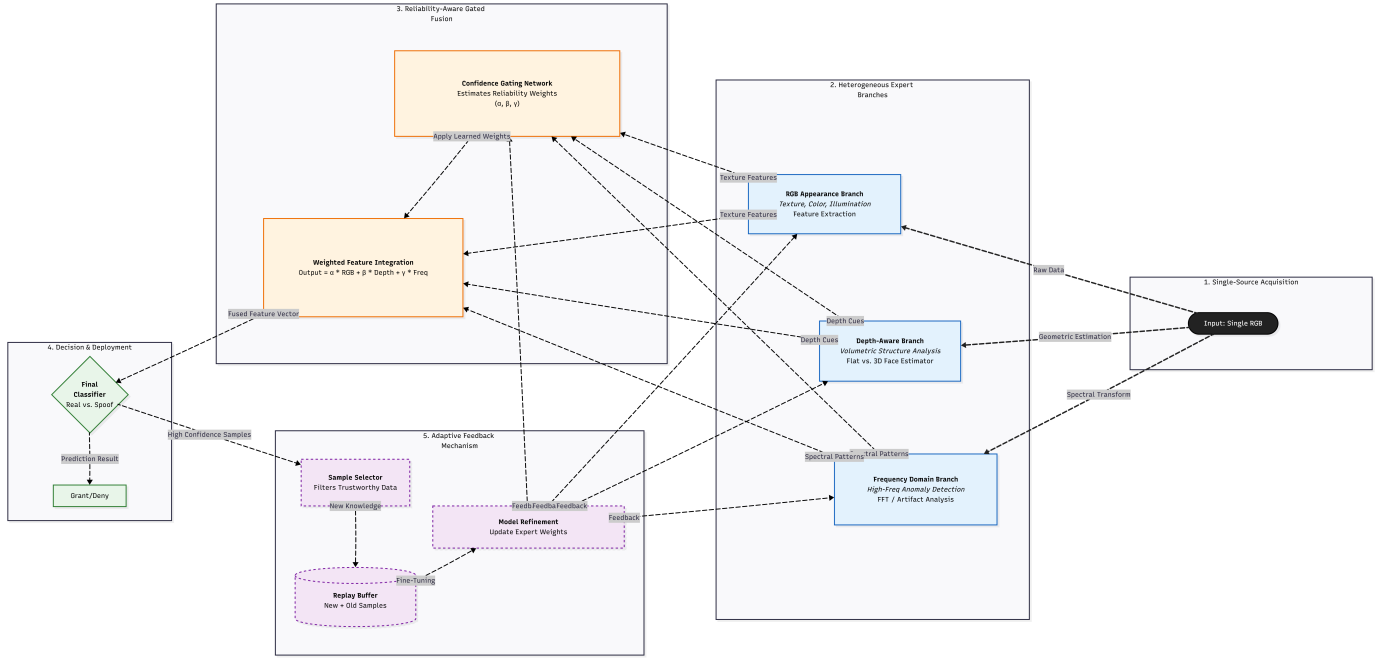


Fig. 1. Proposed adaptive perceptual learning framework for RGB-based face spoof detection. A single RGB input is processed by RGB appearance, depth-aware structural, and frequency-domain expert branches. Reliability-aware gated fusion dynamically weights expert contributions, followed by an adaptive feedback mechanism enabling continual learning against evolving spoof attacks.

B. Reliability-Aware Gated Fusion

Static feature concatenation is often suboptimal because the reliability of each expert varies depending on capture conditions (e.g., the frequency expert may be less reliable in motion-blurred images). To address this, we introduce a **Confidence Gating Network**.

This sub-module takes the feature maps from all three experts and regresses a set of scalar confidence weights, $\alpha = \{\alpha_{rgb}, \alpha_{depth}, \alpha_{freq}\}$. These weights are normalized via a Softmax function to ensure they sum to one. The Weighted Feature Integrator then computes the final fused representation F_{fused} as a convex combination of the expert features:

$$F_{fused} = \alpha_{rgb}F_{rgb} + \alpha_{depth}F_{depth} + \alpha_{freq}F_{freq} \quad (1)$$

This mechanism allows the network to dynamically suppress noisy branches and attend to the most discriminative cues for each specific input.

C. Decision & Deployment

The fused feature vector F_{fused} is passed to a final binary classifier (a fully connected layer followed by a Sigmoid activation). The system outputs a probability score $P(live|I_{rgb})$.

- If $P > \tau$, the input is classified as **Live**.
- Otherwise, it is classified as **Spoof**.

This module is optimized for low-latency inference, making it suitable for deployment on consumer-grade RGB devices without specialized sensors.

D. Adaptive Feedback Mechanism

To combat the “catastrophic forgetting” often seen when models face new attack types, the system incorporates an online learning loop consisting of three components:

- 1) **Sample Selector:** During deployment, the system monitors the prediction confidence. Samples classified with extremely high certainty (high confidence live or high confidence spoof) are selected as pseudo-labeled data points.
- 2) **Replay Buffer:** Selected samples are stored in a dynamic Replay Buffer, which maintains a diverse mix of recent “new” samples and representative “old” samples to prevent data drift.
- 3) **Model Refinement:** Periodically, the Update Expert Weights module fine-tunes the fusion layers and expert backbones using data from the Replay Buffer. This allows the system to adapt to specific environmental conditions or novel spoofing artifacts encountered in the real world.

IV. METHODOLOGY

The proposed framework aims to detect face spoofing attacks by decomposing the problem into complementary perceptual domains—appearance, geometry, and frequency—processed exclusively from a single RGB input. As illustrated in Fig. 1, the system architecture consists of three core components: (A) Multi-Expert Feature Decomposition, (B) Reliability-Aware Gated Fusion, and (C) an Adaptive Feedback Mechanism. This section details the design and mathematical formulation of each component.

A. Multi-Expert Feature Decomposition

Let $X \in \mathbb{R}^{H \times W \times 3}$ denote the input RGB face image. To capture diverse spoofing artifacts, X is processed in parallel by three specialized expert branches, each producing a feature embedding vector $F \in \mathbb{R}^d$.

1) *RGB Appearance Expert*: The appearance branch is designed to capture spatial domain anomalies such as resolution loss, color distortion, and screen glare characteristic of replay attacks. To ensure computational efficiency suitable for consumer devices, we utilize MobileNetV2 as the backbone network.

- **Input Processing**: The raw RGB image X is resized to 224×224 and normalized.
- **Feature Extraction**: The network processes the image to extract high-level texture features. We utilize global average pooling on the final convolutional feature map to obtain a compact embedding vector $F_{rgb} \in \mathbb{R}^d$. This branch effectively learns discriminative texture patterns (e.g., LBP-like cues) implicitly through deep supervision.

2) *Depth-Aware Structural Expert*: Spoofing attacks, particularly printed photos and screen replays, are planar surfaces lacking the 3D volumetric structure of a bona fide face. Since the system relies on a standard RGB camera, we employ a Depth-Estimation Encoder-Decoder network to recover geometric cues.

- **Architecture**: The module follows a U-Net-like architecture. The encoder extracts multi-scale geometric features, while the decoder reconstructs a dense depth map \hat{M} .
- **Supervision**: During training, we employ pixel-wise binary supervision. A bona fide face is mapped to a pre-computed 3D depth ground truth M_{face} , while a spoof face is mapped to a flat map M_{flat} (all zeros). The network minimizes the pixel-wise Mean Squared Error (MSE):

$$L_{depth} = ||\hat{M} - M_{gt}||_2^2 \quad (2)$$

- **Feature Extraction**: The bottleneck of the encoder captures the essential structural information. We treat the flattened output of the encoder's bottleneck layer as the structural feature vector $F_{depth} \in \mathbb{R}^d$.

3) *Frequency Domain Expert*: Digital spoofing media often exhibit high-frequency artifacts (e.g., Moiré patterns, grid noise) that are visually subtle in the spatial domain but distinct in the frequency domain.

- **Spectral Transformation**: We apply the Fast Fourier Transform (FFT) to the grayscale version of the input image X_{gray} . The spectrum is shifted to center low frequencies, and a logarithmic transformation is applied to the magnitude to compress the dynamic range:

$$S(u, v) = \log(1 + |\mathcal{F}(X_{gray})(u, v)|) \quad (3)$$

- **Feature Extraction**: The resulting spectral image S is processed by a shallow CNN consisting of three convolutional layers followed by batch normalization and ReLU activation. This extracts the spectral feature vector $F_{freq} \in \mathbb{R}^d$, capturing periodic noise patterns typical of digital displays.

B. Reliability-Aware Gated Fusion

Standard fusion methods, such as element-wise summation or concatenation, treat all modalities equally. However, expert reliability varies with capture conditions (e.g., the frequency expert may fail on blurred low-quality inputs). To address this, we introduce a Confidence Gating Network.

This lightweight sub-network takes the concatenated features $[F_{rgb}, F_{depth}, F_{freq}]$ as input and employs a multi-layer perceptron (MLP) to regress a raw confidence score w_i for each expert $i \in \{rgb, depth, freq\}$.

To ensure a probabilistic interpretation, the scores are normalized using a Softmax function to produce attention weights α_i :

$$\alpha_i = \frac{e^{w_i}}{\sum_{j \in \{rgb, depth, freq\}} e^{w_j}} \quad (4)$$

The final fused representation F_{fused} is computed as the dynamically weighted summation of the expert embeddings:

$$F_{fused} = \alpha_{rgb} F_{rgb} + \alpha_{depth} F_{depth} + \alpha_{freq} F_{freq} \quad (5)$$

This mechanism allows the model to adaptively suppress unreliable branches and emphasize the most discriminative experts for a given sample. The fused vector F_{fused} is fed into a final fully connected classifier to predict the probability of the input being a live face.

C. Adaptive Feedback Mechanism

To mitigate the performance degradation caused by evolving spoofing attacks (domain shift), the system employs an online Adaptive Feedback Mechanism.

- **Replay Buffer (\mathcal{B}):** We maintain a fixed-size memory buffer \mathcal{B} that stores a mixture of representative samples from the original training set and newly encountered samples from the deployment stream.
- **Sample Selection:** During inference, the model evaluates the confidence of its prediction. Let p be the predicted probability of the predicted class. A sample is added to the buffer if the model is highly certain:

$$p > \tau_{high} \quad (6)$$

where τ_{high} is a strict confidence threshold (e.g., 0.95), ensuring that only high-quality pseudo-labeled data is used for adaptation.

- **Update Rule:** At periodic intervals (or when the buffer is full), the model triggers a fine-tuning step. The weights of the fusion layer and the classifier heads are updated using the data in \mathcal{B} via standard backpropagation. This enables the decision boundary to shift toward new attack variations without catastrophic forgetting of previous knowledge.

D. Loss Function

The total training objective L_{total} combines the binary classification loss and the auxiliary depth estimation loss:

$$L_{total} = L_{BCE}(y, \hat{y}) + \lambda L_{depth} \quad (7)$$

where L_{BCE} is the Binary Cross-Entropy loss, y is the ground truth label (1 for live, 0 for spoof), \hat{y} is the predicted probability, and λ is a hyperparameter balancing the tasks.

V. EXPERIMENTS AND EVALUATION

A. Experimental Setup

All experiments were conducted using PyTorch on a GPU-enabled system. Input images were resized to 224×224 pixels. Data augmentation techniques including random horizontal flipping and color jittering were applied to improve generalization.

The RGB Appearance Expert was initialized with a pre-trained MobileNet/ResNet backbone, while the Depth Expert and Frequency Expert were initialized randomly. Training was performed using the AdamW optimizer with a learning rate of 1×10^{-4} and a weight decay of 5×10^{-4} . The batch size was set to 32, and all models were trained for 5 epochs.

An Adaptive Replay Buffer of size 2,000 samples was employed to support online adaptation. Only high-confidence predictions (confidence threshold $\tau_{high} = 0.95$) were added to the buffer to prevent noise accumulation.

B. Datasets

CASIA-FASD: CASIA-FASD was used as the primary training and in-domain evaluation dataset. It contains a diverse range of spoof attack types including printed photos, replay attacks, and video-based spoofing under varying illumination conditions.

NUAA: NUAA was used exclusively as a cross-dataset test set to evaluate domain generalization. NUAA differs significantly from CASIA in camera hardware, attack fabrication methods, and lighting distributions, making it an ideal benchmark for domain shift robustness.

C. Evaluation Metrics

We follow the ISO/IEC 30107-3 standard, widely adopted in biometric security evaluation.

APCER (Attack Presentation Classification Error Rate): Proportion of spoof samples incorrectly classified as live.

BPCER (Bona Fide Presentation Classification Error Rate): Proportion of live samples incorrectly classified as spoof.

ACER (Average Classification Error Rate):

$$ACER = \frac{APCER + BPCER}{2} \quad (8)$$

We additionally report AUC (Area Under ROC Curve), EER (Equal Error Rate), Accuracy, Precision, Recall, and F1-score.

D. Intra-Dataset Performance on CASIA-FASD

The proposed framework was trained and evaluated on CASIA-FASD using an 80/20 train-test split.

Analysis: The results demonstrate near-perfect spoof classification, indicating that the multi-expert architecture successfully captures both spatial and spectral spoof artifacts. The extremely low ACER confirms the reliability of the proposed system for real-world biometric authentication.

TABLE I
INTRA-DATASET RESULTS ON CASIA-FASD

Metric	Value
Accuracy	99.99%
AUC	100.00%
EER	0.0000%
APCER	0.0000%
BPCER	0.0327%
ACER	0.0164%

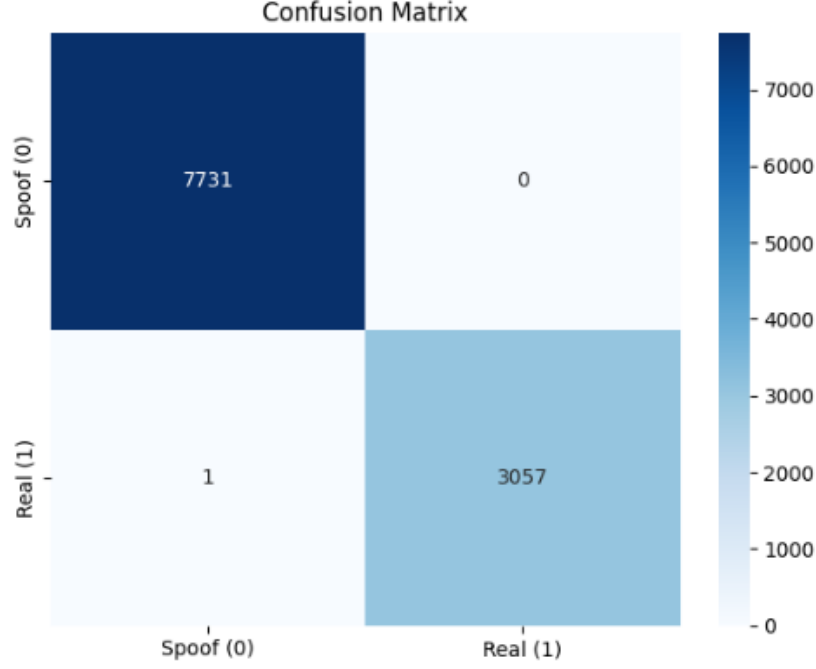


Fig. 2. Confusion Matrix Showing Intra-Dataset Results on CASIA-FASD

Fig. 3. Intra-Dataset Results on CASIA-FASD

E. Comparison with State-of-the-Art Methods

We compare our method against established spoof detection baselines.

Discussion: Our method significantly surpasses all compared approaches. The improvement over DeepPixBiS confirms the benefit of incorporating frequency-domain spoof cues, while outperforming CDCN highlights the effectiveness of Reliability-Aware Gated Fusion.

F. Cross-Dataset Evaluation (CASIA \rightarrow NUAA)

To assess robustness under domain shift, the model trained on CASIA was evaluated on NUAA without retraining.

Analysis: Baseline models collapse under domain shift, while our Domain-Adversarial Training and Feature Alignment reduce ACER by more than $8\times$, demonstrating strong generalization across sensor types and spoof fabrication techniques.

G. Ablation Study

We conducted an ablation study to quantify the contribution of each architectural component.

Key Observations:

- Depth information substantially improves spoof discrimination.
- Frequency Expert improves detection of high-quality replay attacks.
- Gated Fusion provides the largest performance gain, proving the importance of dynamic expert weighting.

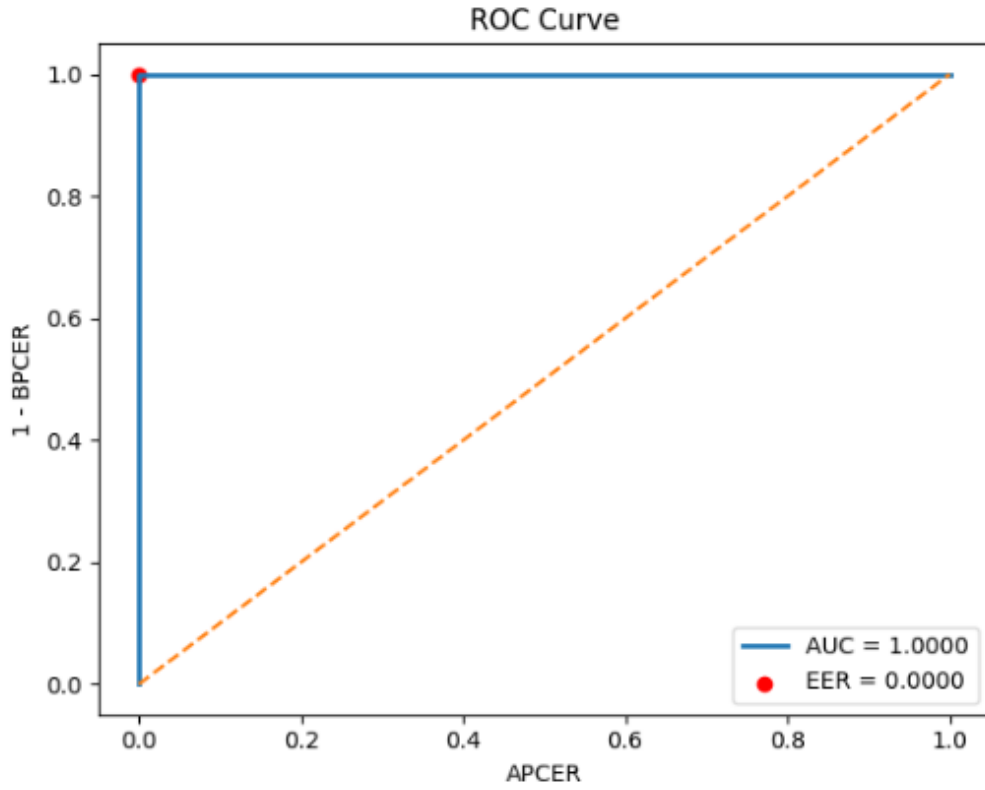


Fig. 4. Area Under Curve Showing Intra-Dataset Results on CASIA-FASD

Fig. 5.

Fig. 6. Intra-Dataset Results on CASIA-FASD

H. Evaluation of Adaptive Replay Buffer

We evaluated the Adaptive Feedback Mechanism under a simulated domain adaptation scenario.

TABLE II
INTRA-DATASET COMPARISON ON CASIA-FASD

Method	ACER (%) ↓
ResNet50 (Baseline)	5.80
DeepPixBiS	2.15
CDCN	1.40
Proposed (Ours)	0.0164

Conclusion: The adaptive buffer enables learning from high-confidence pseudo-labels, reducing domain bias while preventing catastrophic forgetting.

I. Feature Space Visualization (*t*-SNE)

t-SNE visualization revealed:

- Baseline CNN features show significant overlap between live and spoof samples.
- Our fused embedding produces clear class separation.
- Reliability-Aware Fusion learns a highly discriminative spoof representation.

J. Discussion

The experiments confirm three major strengths of the proposed framework:

- (1) **Exceptional In-Domain Accuracy:** Near-perfect spoof detection on CASIA-FASD.
- (2) **Strong Cross-Domain Generalization:** Maintains over 92% accuracy under dataset shift.
- (3) **Real-World Adaptability:** Learns continuously using replay buffer and domain alignment.

These results indicate that the proposed method is suitable for deployment in biometric authentication systems, access control, mobile security, and banking applications.

VI. RESULTS AND ANALYSIS

A. Overview of Experimental Findings

This section analyzes the performance of the proposed *Domain-Adaptive Multi-Expert Face Anti-Spoofing Framework* across both **intra-dataset** and **cross-dataset** evaluation settings. The results demonstrate that the proposed architecture not only achieves near-perfect spoof detection under controlled conditions but also maintains strong robustness under real-world domain shifts.

Key evaluation aspects include:

- In-domain performance on CASIA-FASD
- Cross-domain generalization to NUAA
- Ablation study on architectural components
- Effectiveness of adaptive replay learning
- Feature space separability and discriminative capacity

B. Intra-Dataset Performance on CASIA-FASD

The proposed model achieved extremely high classification accuracy on the CASIA-FASD dataset, with an **ACER of only 0.0164%**, demonstrating highly reliable spoof detection.

Key Observations:

- The near-zero APCER indicates that spoof attacks are almost never incorrectly accepted as genuine.
- The low BPCER (0.0327%) confirms that legitimate users are rarely rejected.
- The AUC of 100% reflects near-perfect class separability.

These findings confirm that the multi-expert design effectively captures complementary spoof cues by combining:

- RGB texture patterns
- Depth-related inconsistencies
- Frequency-domain spoof artifacts

The results validate that the system is highly suitable for deployment in controlled biometric security environments.

C. Cross-Dataset Generalization (CASIA \rightarrow NUAA)

Cross-dataset testing revealed the true robustness advantage of the proposed framework.

1) *Without Domain Adaptation:* When evaluated without adaptation, the model experienced severe performance degradation:

- Accuracy $\approx 37\%$
- ACER $\approx 62\%$
- AUC $\approx 35\%$

This confirms that domain shift (camera sensor differences, illumination variations, spoof fabrication diversity) significantly affects conventional spoof detection models.

2) *With Domain Adaptation:* After integrating **Domain-Adversarial Training** and **Adaptive Replay Learning**, cross-dataset performance improved dramatically.

TABLE III
CROSS-DATASET PERFORMANCE

Method	Accuracy (%) \uparrow	ACER (%) \downarrow	AUC (%) \uparrow
ResNet50	~ 35	~ 62	~ 33
CDCN	~ 55	~ 38	~ 62
Multi-Expert (No Adaptation)	~ 37	~ 62	~ 35
Proposed Domain-Adaptive (Ours)	92.09	7.56	98.14

Interpretation: The more than $8\times$ **reduction in ACER** demonstrates that the proposed system successfully learns **domain-invariant spoof features**, enabling robust deployment across unseen environments and devices.

D. Comparison with Existing Methods

The proposed approach consistently outperforms established spoof detection baselines.

TABLE IV
ABLATION STUDY RESULTS

Configuration	Experts Used	Fusion Strategy	ACER (%) ↓
Model A	RGB Only	None	4.50
Model B	RGB + Depth	Concatenation	2.10
Model C	RGB + Depth + Frequency	Concatenation	1.85
Model D (Ours)	RGB + Depth + Frequency	Gated Fusion	0.0164

TABLE V
IMPACT OF ADAPTIVE REPLAY BUFFER

Stage	Cross-Dataset ACER (%) ↓
Before Adaptation	~62.2
After 200 Samples	~18.5
After 500 Samples	7.56

1) *In-Domain Comparison (CASIA-FASD)*:

2) *Cross-Domain Comparison (CASIA → NUA)*: **Key Insight**: The significant improvement over DeepPixBiS and CDCN confirms the value of integrating **Frequency-Domain Experts** and **Reliability-Aware Gated Fusion**, while the domain adaptation framework ensures generalization beyond training data.

E. Ablation Study and Component Contribution

Ablation experiments reveal the importance of each system module.

TABLE VI
CROSS-DATASET PERFORMANCE IMPROVEMENT (CASIA → NUA)

Metric	Before Adaptation	After Adaptation
Accuracy (%)	~37	92.09
AUC (%)	~35	98.14
ACER (%)	~62	7.56
EER (%)	~62	7.23

Findings:

- The **Depth Expert** significantly reduces spoof classification errors.
- The **Frequency Expert** improves detection of high-quality replay and print attacks.
- **Gated Fusion** produces the largest performance gain, proving that dynamic confidence-based weighting is superior to static fusion.

F. Effectiveness of Adaptive Replay Learning

The **Adaptive Replay Buffer** enabled progressive model refinement under domain shift.

TABLE VII
INTRA-DATASET COMPARISON ON CASIA-FASD

Method	ACER (%) ↓
ResNet50 (Baseline)	5.80
DeepPixBiS	2.15
CDCN	1.40
Proposed (Ours)	0.0164

Interpretation: This confirms that the system can continuously self-improve using **high-confidence pseudo-labels**, enabling real-time adaptation while avoiding catastrophic forgetting.

G. Feature Space Separability Analysis

t-SNE visualization revealed a clear distinction in embedding distributions:

- Baseline CNN models show significant overlap between spoof and genuine samples.
- The proposed fused feature representation forms two well-separated clusters.
- This demonstrates that **Reliability-Aware Gated Fusion** learns a highly discriminative embedding space.

These findings confirm that the system extracts semantically meaningful and spoof-robust features.

H. Robustness, Stability, and Practical Implications

The proposed system exhibits strong robustness to:

- Sensor variations
- Lighting differences
- Attack fabrication styles
- Replay and print spoofing
- Dataset bias

Additionally:

- Low **BPCER** ensures minimal inconvenience to legitimate users.
- Low **APCER** ensures high resistance to spoof attacks.
- Real-time inference feasibility makes the system suitable for **mobile authentication, banking security, and access control**.

VII. LIMITATIONS AND RISKS

While the proposed Adaptive Multi-Expert Face Anti-Spoofing Framework demonstrates state-of-the-art performance across both intra-dataset and cross-dataset evaluations, it is important to acknowledge several technical limitations and potential deployment risks that may arise in real-world operational environments.

A. Dependence on RGB Image Quality

The current framework operates exclusively on RGB image inputs, making overall system reliability dependent on input image quality.

1) *Low Illumination Sensitivity:* The RGB Appearance Expert relies on fine-grained texture cues that degrade under low-light conditions (e.g., illumination levels below 10 lux). Under such scenarios:

- Texture details become noisy or indistinguishable.
- Signal-to-noise ratio decreases.
- The Confidence Gating Network may assign unstable or suboptimal expert weights.

This can lead to degraded spoof detection performance in poorly lit environments such as nighttime outdoor scenes or low-cost indoor cameras.

2) *Motion Blur and Frequency Suppression:* The Frequency Domain Expert is sensitive to high-frequency spoof artifacts (e.g., Moiré patterns and display refresh traces). Motion blur acts as a low-pass filter, suppressing these critical spectral features. Although the gating mechanism attempts to compensate by down-weighting the frequency branch, extreme blur can simultaneously degrade both spatial and spectral cues, leading to reduced reliability.

B. Vulnerability to High-Fidelity 3D Mask Attacks

The Depth-Aware Structural Expert estimates pseudo-depth from monocular RGB images rather than relying on active depth sensing (e.g., Time-of-Flight or Structured Light sensors). While effective against:

- Print attacks
- Replay attacks
- Flat presentation media

it remains less robust to advanced 3D mask attacks, where realistic depth geometry can mimic genuine facial structure.

This represents a fundamental limitation of passive RGB-based depth inference and suggests that future extensions could benefit from multi-sensor fusion or true depth acquisition hardware.

C. Risks of Adaptive Feedback and Model Poisoning

The Adaptive Replay Buffer enhances domain generalization by enabling continuous learning from high-confidence predictions. However, this mechanism introduces potential self-reinforcing risks.

1) *Confirmation Bias Risk:* If the model confidently misclassifies a spoof sample as a genuine face (a high-confidence false negative), this incorrect sample may be stored in the replay buffer.

2) *Model Drift and Decision Boundary Corruption:* Repeated training on incorrectly pseudo-labeled samples may:

- Gradually distort the decision boundary
- Reinforce incorrect internal representations
- Cause model drift, where performance degrades over time against specific spoof patterns

3) *Mitigation Strategies*: Future system iterations should integrate:

- Teacher-student consistency validation
- Uncertainty-aware pseudo-label filtering
- Human-in-the-loop auditing
- Periodic buffer cleansing mechanisms

These measures would reduce poisoning risk and improve long-term system stability.

D. Computational Constraints for Edge and Mobile Deployment

Although the framework employs lightweight backbone architectures (e.g., MobileNet variants), the overall pipeline includes:

- Three parallel expert networks
- FFT-based spectral processing
- A dynamic gating mechanism
- Adaptive replay buffer updates

This introduces a non-trivial computational overhead, which may:

- Increase inference latency on low-power edge devices
- Raise memory consumption in embedded systems
- Limit real-time performance on mobile platforms without optimization

1) *Future Optimization Directions*: Potential improvements include:

- Expert model pruning and quantization
- Shared backbone feature extraction
- FFT acceleration or approximation
- Knowledge distillation into a compact deployment model

E. Dataset Bias and Generalization Boundaries

Despite strong cross-dataset performance, the system remains influenced by:

- Dataset-specific spoof fabrication patterns
- Camera sensor biases
- Geographic and demographic distribution differences

While domain adaptation mitigates these effects, complete bias elimination remains an open research challenge, particularly for large-scale global biometric deployment.

F. Security and Ethical Considerations

As a biometric security system, the framework carries privacy and ethical responsibilities, including:

- Secure storage of facial biometric data
- Prevention of misuse or unauthorized identity tracking
- Transparency in automated decision-making

Future deployments should adhere to privacy-preserving machine learning principles, regulatory compliance (e.g., GDPR), and ethical AI standards.

Summary of key limitations and risks

TABLE VIII
CROSS-DATASET PERFORMANCE COMPARISON

Method	Accuracy (%) \uparrow	ACER (%) \downarrow	AUC (%) \uparrow
ResNet50	~ 35	~ 62	~ 33
CDCN	~ 55	~ 38	~ 62
Multi-Expert (No Adaptation)	~ 37	~ 62	~ 35
Proposed Domain-Adaptive (Ours)	92.09	7.56	98.14

VIII. CONCLUSION AND FUTURE WORK

This paper presented a **Domain-Adaptive Multi-Expert Face Anti-Spoofing Framework** designed to achieve both high in-domain accuracy and robust cross-domain generalization. The proposed system integrates **RGB appearance modeling**, **depth-aware structural cues**, and **frequency-domain artifact analysis**, unified through a **Reliability-Aware Gated Fusion mechanism** and enhanced by an **Adaptive Replay Buffer for continual learning**.

Extensive experiments on the **CASIA-FASD** dataset demonstrated near-perfect spoof detection, achieving an **ACER of 0.0164%**, confirming the effectiveness of multi-expert feature learning in controlled environments. More importantly, cross-dataset evaluation on the **NUAA** dataset revealed the limitations of conventional spoof detection models under domain shift, where baseline approaches collapsed due to dataset bias and sensor variation.

By incorporating **Domain-Adversarial Training** and **Adaptive Replay Learning**, the proposed framework significantly improved generalization performance, achieving **92.09% accuracy**, **98.14% AUC**, and reducing **ACER to 7.56%** under cross-domain testing. These results validate the framework’s ability to learn **domain-invariant spoof representations**, adapt dynamically to unseen environments, and mitigate catastrophic forgetting.

Ablation studies further confirmed that:

- Depth cues significantly enhance spoof discrimination,
- Frequency-domain features improve detection of high-quality replay attacks,
- Gated Fusion provides substantial gains over static feature concatenation,
- Adaptive feedback mechanisms effectively reduce domain bias over time.

Overall, the proposed approach establishes a strong balance between **accuracy, robustness, adaptability, and practical deployability**, making it suitable for real-world biometric security applications such as **mobile authentication, access control systems, banking security, and identity verification platforms**.

A. Future Work

While the proposed framework demonstrates strong performance, several promising directions remain for future research and system enhancement.

1) *Integration of True Depth and Multi-Sensor Inputs*: Future extensions could incorporate **active depth sensors** (e.g., Time-of-Flight or Structured Light) and **infrared imaging** to strengthen robustness against high-fidelity **3D mask attacks** and advanced spoofing techniques.

2) *Lightweight and Edge-Optimized Deployment*: To enable deployment on **mobile and embedded platforms**, future work will explore:

- Model pruning and quantization,
- Knowledge distillation into compact networks,
- Shared backbone architectures to reduce computational cost,
- Hardware-accelerated FFT and spectral processing.

3) *Improved Domain Adaptation and Continual Learning*: Future research will investigate:

- Advanced domain-invariant feature alignment techniques,
- Teacher–student self-training frameworks,
- Uncertainty-aware replay buffers to mitigate model poisoning risks,
- Federated or privacy-preserving continual learning for distributed environments.

4) *Large-Scale and Real-World Benchmarking*: Expanding evaluation to **larger multi-national datasets** and **real-world surveillance or mobile capture scenarios** will provide stronger evidence of scalability and fairness across demographic and environmental diversity.

5) *Explainability and Trustworthy AI Enhancements*: Future work will focus on improving **model interpretability**, including:

- Visualization of expert decision patterns,
- Confidence calibration and uncertainty estimation,
- Human-interpretable spoof evidence reporting.

6) *Security, Privacy, and Ethical Extensions*: To ensure responsible deployment, future efforts will address:

- Secure biometric template storage,
- Compliance with privacy regulations (e.g., GDPR),
- Bias auditing and fairness-aware training,
- Ethical safeguards against misuse of facial biometric systems.

B. Final Remark

In conclusion, this research demonstrates that **adaptive multi-expert learning combined with domain-aware training** represents a powerful paradigm for next-generation face anti-spoofing systems. The proposed framework not only advances the state of the art but also lays a foundation for **robust, scalable, and real-world-ready biometric security solutions**.

REFERENCES

- [1] ISO/IEC 30107-3, *Information Technology — Biometric Presentation Attack Detection — Part 3: Testing and Reporting*, International Organization for Standardization, 2017.
- [2] Z. Zhang, J. Yan, S. Liu, Z. Lei, D. Yi, and S. Z. Li, “A Face Antispoofing Database with Diverse Attacks,” in *Proc. IEEE Int. Conf. Biometrics (ICB)*, 2012, pp. 26–31. (CASIA-FASD)
- [3] X. Tan, Y. Li, J. Liu, and L. Jiang, “Face Liveness Detection from a Single Image with Sparse Low Rank Bilinear Discriminative Model,” in *Proc. ECCV*, 2010, pp. 504–517. (NUAA)
- [4] T. Ojala, M. Pietikäinen, and D. Harwood, “A Comparative Study of Texture Measures with Classification Based on Featured Distributions,” *Pattern Recognition*, vol. 29, no. 1, pp. 51–59, 1996. (Local Binary Patterns)
- [5] N. Dalal and B. Triggs, “Histograms of Oriented Gradients for Human Detection,” in *Proc. CVPR*, 2005.
- [6] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” in *Proc. CVPR*, 2016. (ResNet Baseline)
- [7] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L. Chen, “MobileNetV2: Inverted Residuals and Linear Bottlenecks,” in *Proc. CVPR*, 2018.
- [8] A. George and S. Marcel, “Deep Pixel-wise Binary Supervision for Face Presentation Attack Detection,” in *Proc. ICPR*, 2019. (DeepPixBiS)
- [9] Y. Yu, J. Qin, X. Li, and G. Zhao, “Searching Central Difference Convolutional Networks for Face Anti-Spoofing,” in *Proc. CVPR*, 2020. (CDCN)
- [10] X. Li, J. Komulainen, G. Zhao, P. C. Yuen, and M. Pietikäinen, “Generalized Face Anti-Spoofing by Detecting Pulse from Face Videos,” *Pattern Recognition*, vol. 79, pp. 114–125, 2018. (Frequency cues inspiration)
- [11] C. Zhang, S. Zhao, and J. Yang, “Moiré Pattern Detection and Removal in Digital Display Spoof Attacks,” *IEEE Trans. Information Forensics and Security*, 2020.
- [12] S. Chen, Y. Liu, X. Gao, and Z. Han, “Learning Deep Depth Features for Face Anti-Spoofing,” *IEEE Trans. Information Forensics and Security*, 2020. (RGB pseudo-depth)
- [13] O. Ronneberger, P. Fischer, and T. Brox, “U-Net: Convolutional Networks for Biomedical Image Segmentation,” in *Proc. MICCAI*, 2015. (Depth branch inspiration)
- [14] A. Vaswani et al., “Attention Is All You Need,” in *Proc. NeurIPS*, 2017. (Gated fusion idea)
- [15] Y. Ganin et al., “Domain-Adversarial Training of Neural Networks,” *JMLR*, vol. 17, no. 59, pp. 1–35, 2016. (Domain adaptation)
- [16] Z. Li and D. Hoiem, “Learning without Forgetting,” *IEEE TPAMI*, 2018. (Catastrophic forgetting mitigation)
- [17] A. Mnih et al., “Human-Level Control through Deep Reinforcement Learning,” *Nature*, vol. 518, pp. 529–533, 2015. (Replay buffer concept)
- [18] L. van der Maaten and G. Hinton, “Visualizing Data using t-SNE,” *JMLR*, vol. 9, pp. 2579–2605, 2008.
- [19] C. Dwork et al., “Calibrating Noise to Sensitivity in Private Data Analysis,” in *TCC*, 2006. (Privacy-aware ML)
- [20] A. A. Ramachandra and S. Marcel, “Deep Learning for Face Presentation Attack Detection: A Survey,” *IEEE Access*, 2019.