



PDF Download
3746027.3754939.pdf
14 January 2026
Total Citations: 0
Total Downloads: 78

 Latest updates: <https://dl.acm.org/doi/10.1145/3746027.3754939>

RESEARCH-ARTICLE

InstructFLIP: Exploring Unified Vision-Language Model for Face Anti-spoofing

KUN HSIANG LIN, National Taiwan University, Taipei, Taiwan

YUWEN TSENG, National Taiwan University, Taipei, Taiwan

KANGYANG HUANG, National Taiwan University, Taipei, Taiwan

JHIH CIANG WU, National Taiwan Normal University, Taipei, Taiwan

WENHUANG CHENG, National Taiwan University, Taipei, Taiwan

Open Access Support provided by:

National Taiwan University

National Taiwan Normal University

Published: 27 October 2025

[Citation in BibTeX format](#)

MM '25: The 33rd ACM International
Conference on Multimedia
October 27 - 31, 2025
Dublin, Ireland

Conference Sponsors:
SIGMM

InstructFLIP: Exploring Unified Vision-Language Model for Face Anti-spoofing

Kun-Hsiang Lin
National Taiwan University
Taipei, Taiwan
jacklin@cmlab.csie.ntu.edu.tw

Yu-Wen Tseng
National Taiwan University
Taipei, Taiwan
d12922018@csie.ntu.edu.tw

Kang-Yang Huang
National Taiwan University
Taipei, Taiwan
huangkangyang@cmlab.csie.ntu.edu.tw

Jhih-Ciang Wu
National Taiwan Normal University
Taipei, Taiwan
jcwu@csie.ntnu.edu.tw

Wen-Huang Cheng*
National Taiwan University
Taipei, Taiwan
wenhuang@csie.ntu.edu.tw

Abstract

Face anti-spoofing (FAS) aims to construct a robust system that can withstand diverse attacks. While recent efforts have concentrated mainly on cross-domain generalization, two significant challenges persist: *limited semantic understanding of attack types* and *training redundancy across domains*. We address the first by integrating vision-language models (VLMs) to enhance the perception of visual input. For the second challenge, we employ a meta-domain strategy to learn a unified model that generalizes well across multiple domains. Our proposed *InstructFLIP* is a novel instruction-tuned framework that leverages VLMs to enhance generalization via textual guidance trained solely on a single domain. At its core, InstructFLIP explicitly decouples instructions into content and style components, where content-based instructions focus on the essential semantics of spoofing, and style-based instructions consider variations related to the environment and camera characteristics. Extensive experiments demonstrate the effectiveness of InstructFLIP by outperforming SOTA models in accuracy and substantially reducing training redundancy across diverse domains in FAS. The project website is available at <https://kunkunlin1221.github.io/InstructFLIP>.

CCS Concepts

• Computing methodologies → Biometrics.

Keywords

Face Anti-spoofing, Vision-language models, Unified model

ACM Reference Format:

Kun-Hsiang Lin, Yu-Wen Tseng, Kang-Yang Huang, Jhih-Ciang Wu, and Wen-Huang Cheng. 2025. InstructFLIP: Exploring Unified Vision-Language Model

*Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '25, Dublin, Ireland

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-2035-2/2025/10
<https://doi.org/10.1145/3746027.3754939>

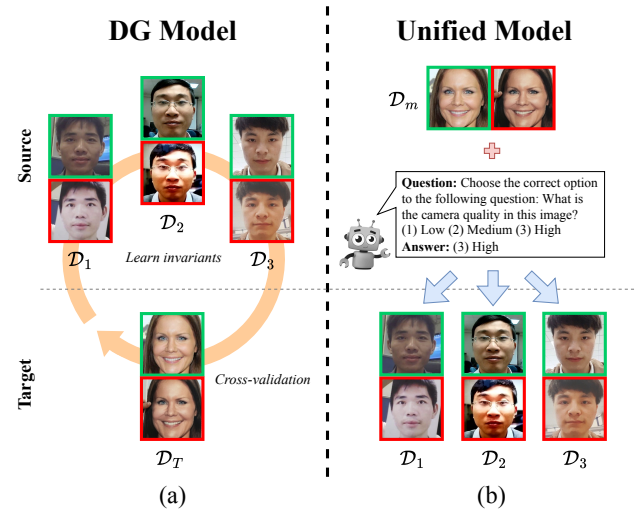


Figure 1: Comparison between DG-based and unified model. (a) Conventional DG-based approaches using leave-one-out protocols from multiple domains ($\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3$) to the target one (\mathcal{D}_T) result in significant training redundancy by treating each domain independently during cross-validation. (b) We propose InstructFLIP, a unified model that leverages the meta domain \mathcal{D}_m and instruction-based training to jointly learn domain-invariant content and style features, allowing efficient generalizations without redundant retraining.

for Face Anti-spoofing. In *Proceedings of the 33rd ACM International Conference on Multimedia (MM '25)*, October 27–31, 2025, Dublin, Ireland. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3746027.3754939>

1 Introduction

Modern face recognition systems have become pervasive, from unlocking smartphones to accessing secure facilities, due to advancements in deep learning and their non-contact, user-friendly nature. However, these systems are susceptible to presentation attacks, where adversaries attempt to bypass security mechanisms using tools such as printed photos, replayed videos, and other advanced deception techniques. Consequently, FAS plays a crucial role in safeguarding these systems against such vulnerabilities. While

existing FAS methods exhibit strong performance in intra-dataset evaluations, where training and testing data share the same domain, their effectiveness significantly declines in cross-dataset scenarios due to domain distribution shifts [41].

To ensure the reliability of FAS in real-world applications, considerable research has concentrated on domain generalization (DG) techniques [7, 8, 17, 32, 34, 37, 53]. DG aims to facilitate models in learning generalized feature representations that effectively distinguish between genuine and spoofed inputs across diverse and unpredictable environments, thereby enhancing both robustness and adaptability. Such approaches are indispensable for developing systems capable of maintaining high performance even in the presence of novel, previously unseen attack scenarios. However, DG requires extensive cross-validation across datasets, resulting in training redundancy due to repeated retraining. In Figure 1(a), conventional DG methods exacerbate this issue by using leave-one-out protocols, where models are trained on multiple domains like (\mathcal{D}_1 , \mathcal{D}_2 , \mathcal{D}_3) and evaluated on a target domain (\mathcal{D}_T), treating each domain independently. This inefficiency limits the practicality of DG for large-scale deployment as they struggle to produce a unified, generalizable model. Additionally, DG often lacks transparency, providing limited insight into their decision-making processes.

Recent advancements in VLMs, such as CLIP [28], aim to improve generalization through cross-modal learning. For example, FLIP [30] leverages the pretrained capabilities of CLIP to address the FAS task, while CFPL [23] employs BLIP2's Q-Former [18] to bridge the modality gap between images and pretrained text encoders. However, these VLM-based methods have not fully exploited the rich semantic information inherent in language, limiting their effectiveness in complex scenarios requiring meticulous understanding.

In summary, two fundamental limitations remain unaddressed in existing DG and VLM approaches: (1) *limited semantic understanding of facial attacks*, as many methods struggle to accurately understand facial attacks due to interference from factors such as environmental conditions and camera quality. While most rely on implicit patterns to propose solutions, they lack intuition, making reliable decision-making challenging; and (2) *training redundancy across domains*, where traditional leave-one-out protocols require models to be trained and validated independently on multiple domains, leading to excessive retraining. To address these challenges, we propose InstructFLIP, a novel VLM-based architecture that applies instruction tuning to explicitly extract and encode high-level semantic priors from textual descriptions, as illustrated in Figure 1(b). Effective FAS necessitates the learning of domain-invariant representations while preserving sensitivity to discriminative, non-spoof-specific attributes. By conditioning the model on structured linguistic instructions, InstructFLIP enhances its capacity to align visual features with task-relevant semantics, thereby improving generalization and robustness across heterogeneous domains.

To more effectively encode spoof-related and non-spoof-related attributes within the VLM framework, we decompose linguistic supervision into two orthogonal components: content and style. The content component encodes spoof-specific semantics such as replay, print, mask attacks, etc., while the style component captures nuisance factors unrelated to spoofing, including image quality, environmental context, and illumination conditions. This structured decomposition facilitates disentangled feature learning and

enhances the model's robustness to domain shifts. Aiming to avoid training redundancy, we introduce a richly labeled meta-domain \mathcal{D}_m designed to synthesize diverse image-instruction pairs for efficient and unified model training. As depicted in Figure 1(b), our approach leverages \mathcal{D}_m in conjunction with instruction tuning to jointly learn domain-invariant content and style representations, eliminating the need for repeated retraining across domains. This paradigm significantly reduces training overhead while improving generalization capability, thereby enhancing the model's scalability for real-world deployment. By consolidating diverse supervisory signals into a single training cycle, InstructFLIP effectively captures multiple task-relevant factors, enabling adaptive and generalized FAS. The main contributions of this paper are as follows:

- This paper proposes InstructFLIP, a novel instruction-tuned VLM framework for FAS, which integrates textual supervision to enhance semantic understanding of spoofing cues.
- We design a content-style decoupling mechanism that explicitly separates spoof-related (content) and spoof-irrelevant (style) information, improving generalization to unseen domains.
- We introduce a meta-domain learning strategy to eliminate training redundancy in cross-domain settings by utilizing diverse image-instruction pairs sampled from a structured meta-domain.
- Experimental results demonstrate that InstructFLIP surpasses SOTA methods across multiple FAS benchmarks, effectively capturing spoof-related patterns through language-guided supervision while substantially reducing training overhead, thereby enhancing its applicability in real-world scenarios.

2 Related Work

2.1 Face Anti-Spoofing

FAS is a vital safeguard in face recognition systems aiming to mitigate the risks posed by presentation attacks. Over recent decades, research in this field has significantly evolved, particularly with the advancement of deep learning. The focus has shifted towards leveraging neural networks to enhance detection accuracy. This includes the development of classification-based models [1, 25, 27, 39, 40] and other auxiliary-learning methods [33, 42], and generative frameworks [31, 36] to improve robustness against spoofing attempts.

The increasing complexity of cross-domain variability, driven by diverse application scenarios and heterogeneous datasets, has led to the development of domain adaptation (DA) techniques [35, 43, 54]. These methods enhance model adaptability by learning shared feature representations across source and target domains. However, DA approaches often require extensive fine-tuning and substantial labeled data, which limits their scalability and practicality for real-world FAS applications. To overcome these challenges, DG techniques [17, 32, 34, 37, 47, 52] have been developed to extract domain-invariant features, enabling models to generalize to unseen domains. Despite their potential, DG methods necessitate extensive cross-validation across multiple datasets, resulting in significant training inefficiency and high costs due to repeated retraining. This constraint hampers their practicality for large-scale deployment, as they not only fail to produce a unified, adaptable model but also often lack the reasoning ability needed for reliable decision-making.

Therefore, we propose InstructFLIP, which eliminates the need for redundant retraining by utilizing a single meta-domain for

Table 1: Options of content-related and style-related prompts.

Type	Question	Options
Content	Spoof type	(1) Real face (2) Photo (3) Poster (4) A4-paper (5) 2D face mask (6) 2D upper-body mask (7) 2D region mask (8) PC screen (9) Pad screen (10) Phone screen (11) 3D mask
Style	Illumination	(1) Normal (2) Strong (3) Back (4) Dark
	Environment	(1) Indoor (2) Outdoor
	Camera quality	(1) Low (2) Medium (3) High

unified model learning. By incorporating instruction tuning with rich linguistic guidance, InstructFLIP achieves improved semantic alignment and robust generalization across diverse domains.

2.2 Visual Language Models

Since the introduction of CLIP [28], VLMs [6, 19] have advanced rapidly, yielding significant breakthroughs across multiple research fields [11, 46, 51]. In the context of FAS, researchers have explored the integration of textual information to facilitate DG. For instance, FLIP [30] pioneered the application of CLIP to address FAS challenges, while extensions such as [3, 12, 13, 21, 23, 24, 44, 55] sought to enhance the generalization capabilities of FAS models through visual language learning. Nevertheless, these approaches have primarily targeted cross-domain issues without fully leveraging the potential of textual understanding to improve model understanding and performance in unseen scenarios.

To address the limitations of existing approaches, this work presents InstructFLIP, a novel instruction-tuned methodology inspired by InstructBLIP [6]. This approach leverages the rich annotations available in the CelebA-Spoof (CS) dataset [49] to construct structured instruction-based image-text pairs, enabling the training of a unified FAS model. By incorporating comprehensive and semantically meaningful instructions, the model learns to generalize effectively across diverse domains while gaining language-level representations of spoofing mechanisms. This not only enhances robustness and improves the accurate identification of spoof types, but also underscores the critical role of reasoning capabilities in VLM for effective FAS solutions.

3 Proposed Method

Learning to perform FAS often presents DG challenges, requiring significant training redundancy to accommodate diverse protocol requirements. To address this, we aim to develop a unified approach capable of effectively adapting to various scenarios. Specifically, we propose **Instruct-tuned Face Anti-spoofing Language-Image Pre-training (InstructFLIP)**, a novel framework fine-tuned on a meta-domain \mathcal{D}_m , which enables DG without the need for repeated retraining, as illustrated in Figure 2. By decoupling the representation into content and style features, our model can adapt to varying domain characteristics by leveraging corresponding interpretable components, promoting a comprehensive visual understanding that enhances generalization to handle domain variations. The learned embeddings from the content and style branches are further fused with the visual encoder via a dedicated fusion module to produce the final prediction, enabling LLM-free inference and supporting

real-time deployment. In the following section, we elaborate on our InstructFLIP in addressing FAS, which facilitates robust generalization across diverse domains.

3.1 Instructive Prompts

To efficiently capture versatile semantic representations for dealing with diverse domains, we restrict our training data solely from \mathcal{D}_m [49], enable to learn a unified model that generalizes effectively. By leveraging the instances from \mathcal{D}_m , we aim to build a model capable of learning domain-invariant features, which facilitates robust performance across varying domain conditions without needing extensive retraining on each specific domain.

Building on insights from previous literature [23, 37, 52], we structure the instructive prompts into two categories, *i.e.*, content and style prompts. The former targets to identify the spoofing types, while the latter focuses on style conditions, such as illumination, environment, and camera quality, as introduced by \mathcal{D}_m . More precisely, we craft instruction-based questions for content prompts, such as “Choose the correct option for the following question: Which spoof type is in this image?” The options provided cover a range of spoofing types along with the real face category, as outlined in Table 1. For style prompts, we use a similar structure to content prompts but concentrate on style attributes. For example: “Choose the correct option for the following question: What is the illumination condition/environment/camera quality in this image?” The options here are tailored to address variations in lighting, background, and camera quality, as indicated by style questions. By exploring \mathcal{D}_m with these prompts, we create a varied set for enriching texture inputs that are crucial while training the unified model. The composition of content and style prompts ensures that the model is equipped with a comprehensive understanding of both the visual characteristics that influence spoofing detection, improving its general performance across various environmental conditions.

3.2 Visual Content and Style Features

Since the instruction prompts are divided into content and style classes, we extend this concept to disentangle features. Given a face image x , we represent $f_c = E(x)$ as the content feature encoded by a feature extractor E and obtained from the highest-level layer. To construct the style feature, we adopt an approach inspired by Adaptive Instance Normalization [16] to assemble a multi-level feature. Unlike conventional method [23] that aggregates style features from each level through averaging, we propose retaining distinct characteristics from each level. Specifically, we concatenate the mean and standard deviation of the features from multiple levels, which can be formulated as

$$f_s = s_1 \oplus s_2 \cdots \oplus s_L, \quad (1)$$

where $s_i = \mu_i \oplus \sigma_i$ stands for a concatenation of the mean and standard deviation in i -level. μ_i and σ_i are obtained by calculating the statistic value of $E_i(x)$.

We clarify the notations by elaborating on the dimensions of f_c and f_s for a clear understanding. The content feature $f_c \in \mathbb{R}^{N \times d}$ is a tensor, where N represents the total number of tokens, including the class and patch tokens, and d is the dimensionality of each

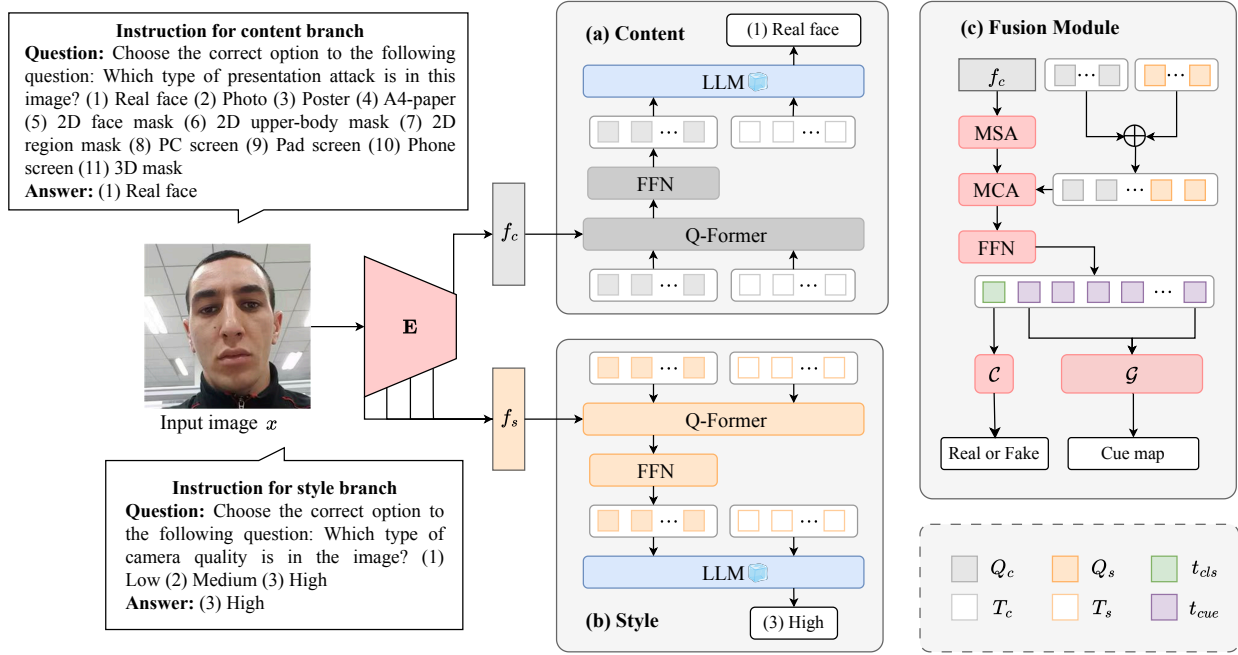


Figure 2: Overview of the proposed InstructFLIP framework for FAS. We first encode the input image x , resulting in the separated content and style features f_c and f_s . The dual-branch architecture presented in (a) and (b) is designed for instruction tuning based on the corresponding expertise in each branch. Eventually, the prediction for determining x is spoofing or not is carried out in (c) via a classifier C according to the fused queries and f_c , coupling with the cue map generated by G .

token embedding. The style feature $f_s \in \mathbb{R}^{2L \times d}$ is a multi-level representation, where L denotes the total number of layers in E .

3.3 Instruction Tuning with Visual Features

We employ a couple of features, *i.e.*, f_c and f_s , by proposing a dual-branch architecture that learns distinct representations for each aspect. The content branch attempts to capture the attributes directly related to attack types, while the style branch gathers contextual information not directly associated with spoofing but is vital for understanding scene variability. As depicted in Figure 2, the content and style branches share a similar architectural structure. Here, we primarily describe the process within the content branch, where the equivalent operations can likewise be applied to the style branch. The instructive content prompt first goes through the standard process, such as tokenization and word embedding, resulting in an instruction representation $T_c \in \mathbb{R}^{j \times d}$. This instruction T_c is combined with the learnable query $Q_c \in \mathbb{R}^{k \times d}$ to compose the input for the Q-Former [18], where j and k represent the length of instruction and the number of learnable queries, respectively. The Q-Former stacks multiple blocks, with each containing multi-head self-attention (MSA) and multi-head cross-attention (MCA), where we symbolize them by ϕ and ψ for simplicity. For each hidden state h , the MSA is defined as

$$\tilde{h} = h + \phi(h), \quad (2)$$

where $h = Q_c \oplus T_c$ is the initial state for the Q-Former, which is subsequently updated. We incorporate the hidden state with

content feature via MCA, which can be formulated as

$$\hat{h} = h + \psi(h, f_c). \quad (3)$$

To streamline the expression, we omit layer normalization typically applied to each hidden state and the content feature before the attention mechanism and disregard the two linear layers with the activation function expanded after the MCA.

After passing through the Q-Former, the query effectively learns to capture critical affinities from the content feature. The processed query is subsequently concatenated with the instruction representation, serving as a soft prompt to frozen LLMs to generate the corresponding prediction p_c , *i.e.*, spoofing type or real face in the content branch.

The objectives in the two branches seek that the model effectively distinguishes between spoof-related and domain-invariant characteristics. For example, the objective in the content branch is optimized using cross-entropy loss, accurately classifying spoof-related attributes by minimizing the discrepancy between prediction and ground truth, which can be formulated as

$$\mathcal{L}_c = - \sum y_c \log(p_c), \quad (4)$$

where y_c and p_c represent the label and prediction, respectively. The style branch is optimized by \mathcal{L}_s similarly to (4).

3.4 Content Feature-assist Query Fusion

To combine the content and style information for spoof detection, we comprise the learned queries from each branch using the attention mechanisms presented in (2) and (3). Precisely, we reuse

Table 2: Unified FAS benchmark results on MCIO and WCS datasets. The CA dataset is used for training, with evaluation conducted on both MCIO and WCS benchmarks. Subtable (a) presents results on MCIO, (b) reports on WCS, and (c) summarizes average performance across all datasets. The best and second-best results are highlighted in red and blue, respectively.

(a) Results on MCIO datasets													
Method	Venue	M			C			I			O		
		HTER↓	AUC↑	TPR↑@ FPR=1%	HTER↓	AUC↑	TPR↑@ FPR=1%	HTER↓	AUC↑	TPR↑@ FPR=1%	HTER↓	AUC↑	TPR↑@ FPR=1%
SSDG-R [17]	CVPR'20	24.92	82.42	16.19	17.84	89.77	21.36	18.60	90.00	47.67	25.68	81.19	22.74
ViT [14]	ECCV'22	12.69	91.42	67.93	7.36	97.09	74.18	20.08	87.04	64.29	25.48	82.18	40.58
SAFAS [32]	CVPR'23	31.82	74.01	3.49	29.26	75.17	4.20	27.57	79.23	15.90	29.84	76.38	13.76
FLIP-MCL [30]	ICCV'23	13.60	91.70	59.29	5.25	98.99	91.11	17.37	91.46	65.18	17.72	89.02	40.02
CFPL [23]	CVPR'24	8.45	95.48	73.31	2.77	99.53	94.36	10.00	95.90	80.13	17.21	89.30	31.58
InstructFLIP	-	5.52	98.12	86.62	1.47	99.79	97.36	9.12	96.17	86.48	14.33	94.79	64.77

(b) Results on WCS datasets										(c) Mean metrics computed over all datasets.			
Method	W			C			S			Method	Avg.		
	HTER↓	AUC↑	TPR↑@ FPR=1%	HTER↓	AUC↑	TPR↑@ FPR=1%	HTER↓	AUC↑	TPR↑@ FPR=1%		HTER↓	AUC↑	TPR↑@ FPR=1%
SSDG-R [17]	41.57	62.88	2.38	19.98	87.87	20.38	52.41	46.87	1.09	SSDG-R [17]	28.71	77.29	18.83
ViT [14]	28.08	78.49	6.89	12.50	94.52	49.71	39.81	64.62	3.65	ViT [14]	20.86	85.05	43.89
SAFAS [32]	41.69	61.11	0.97	18.14	89.06	23.34	55.51	41.29	0.76	SAFAS [32]	33.40	70.89	8.92
FLIP-MCL [30]	26.96	79.96	11.96	7.90	97.54	57.31	42.97	59.89	0.47	FLIP-MCL [30]	18.82	86.94	46.48
CFPL [23]	26.85	80.97	17.42	5.50	98.74	78.72	42.29	60.60	2.37	CFPL [23]	16.15	88.65	53.98
InstructFLIP	19.51	89.90	39.11	8.49	96.81	70.17	30.33	80.16	12.07	InstructFLIP	12.68	93.68	65.23

MSA and MCA operations to integrate the content and style queries along with the content feature f_c , creating a unified representation for spoof detection. The integration process is defined as

$$\hat{Q} = Q + \psi(Q, \tilde{f}_c), \quad (5)$$

where $Q = Q_c \oplus Q_s \in \mathbb{R}^{2k \times d}$ is the concatenated representation of queries from content and style branches. Note that we do not include the style feature f_s in the fusion process, as it is a domain-related representation and could potentially conflict with our goal of achieving robust generalization in the unified model.

The resulting unified representation \hat{Q} from (5) is divided into two segments. The first part contains a single foremost token, denoted as t_{cls} , which serves as the primary representation for the face spoofing classifier C . The second part, t_{cue} , includes the remaining tokens, which capture auxiliary components and are fed into the cue generator \mathcal{G} designed to produce attack hints to enhance the robustness. Given the uncertain nature of attack cues and the unavailable annotations within the dataset, we take inspiration from one-class training paradigms [15]. Specifically, we apply a convolution layer that introduces guided, normalized noise, encouraging the model to generate a uniformly white map for fake samples, which signifies the presence of an attack on the input face. Conversely, the model is trained to produce an entirely blank map for real samples, indicating the absence of spoofing to assist in addressing FAS.

3.5 Objectives

Our proposed InstructFLIP is trained with the loss function that combines multiple components, which can be represented as

$$\mathcal{L} = \lambda_1 \mathcal{L}_c + \lambda_2 \mathcal{L}_s + \lambda_3 \mathcal{L}_{cls} + \lambda_4 \mathcal{L}_{cue}, \quad (6)$$

where λ_i are the hyperparameters for balancing each object. The first two terms, \mathcal{L}_c and \mathcal{L}_s , are introduced in content and style branches for carrying out instruction tuning. The third objective, \mathcal{L}_{cls} optimizes the prediction for the face spoofing classifier C by cross-entropy. The last term, \mathcal{L}_{cue} , is designed for learning the cue map from \mathcal{G} , which is defined as

$$\mathcal{L}_{cue} = \begin{cases} \mathbf{d}^2 / \beta & \text{if } \mathbf{d} < \beta \\ \mathbf{d} - \beta / 2 & \text{otherwise,} \end{cases} \quad (7)$$

where $\mathbf{d} = |p_{cue} - y_{mask}|$ denotes the distance between the predicted cue map and GT binary mask. β is the threshold at which to change between L1 and L2 loss.

4 Experiment

4.1 Experimental Setup

Datasets and Evaluation Protocol. To rigorously evaluate the proposed framework, we adopt a unified evaluation protocol that departs from traditional leave-one-out or one-to-one strategies [14, 23, 30, 32], aiming to train a single model with strong cross-domain generalization while avoiding training redundancy. The CelebA-Spoof (CA) [49] is used for training, offering rich attribute-level annotations that enable the generation of diverse instruction-based image-text pairs. Specifically, eleven fine-grained spoof labels are

used to construct textual supervisions for the content branch, while three imaging conditions, illumination, environment, and camera quality, serve as style-based contexts based on the extensive metadata available in the CA. Generalization performance is evaluated on seven publicly available FAS benchmarks: MSU-MFSD (M)[38], CASIA-FASD (C)[50], Replay-Attack (I)[4], OULU-NPU (O)[2], WMCA (W)[10], CASIA-CeFA (C)[20, 22], and CASIA-SURF (S) [47, 48]. Under this unified setting, the model is trained once and evaluated on all target datasets concurrently, thereby providing a comprehensive evaluation of its generalization performance.

Implementation Details. Face images are aligned by MTCNN [45] and resized to $224 \times 224 \times 3$. We adopt the ViT-B/16 from CLIP [28] as the visual feature extractor E and use FLAN-T5 base [5] as the frozen LLMs. Optimization is performed using the AdamW optimizer [26] in combination with the OneCycleLR learning rate scheduler [29]. The initial learning rate is set to 1×10^{-6} , with a peak learning rate of 5×10^{-6} and a weight decay of 10^{-6} . All models are trained end-to-end for 20 epochs on 1 NVIDIA RTX 4090 GPU with a batch size of 24. The loss weights in InstructFLIP are set as follows: $\lambda_1 = 0.4$, $\lambda_2 = 0.4$, $\lambda_3 = 0.15$, and $\lambda_4 = 0.05$.

Evaluation Metrics. Following prior work [14, 23, 30, 32], we evaluate model performance using three standard metrics: Half Total Error Rate (HTER), Area Under the Receiver Operating Characteristic Curve (AUC), and True Positive Rate (TPR) at a fixed False Positive Rate (FPR). To ensure fair and consistent evaluation, each model is trained five times with different random seeds, and the mean performance across runs is reported for all metrics.

4.2 Unified FAS Performance

Baseline Methods. The primary and most competitive baselines for evaluating the proposed InstructFLIP method include SSDG [17], ViT [14], SAFAS [32], FLIP [30], CFPL [23], BUDoPT [24], Diff-FAS [9], and FGPL [13]. While ViT provides results for both zero-shot and few-shot settings, we report only the zero-shot setting to maintain consistency with our unified protocol, and we exclude the adaptive ViTAF variant specific to the few-shot scenario. Furthermore, due to the unavailability of publicly released codes for BUDoPT, DiffFAS, and FGPL, these methods are excluded from our experimental comparisons. Consequently, we reimplemented SSDG, ViT, SAFAS, FLIP, and CFPL under our unified evaluation framework and presented their results in the subsequent subsections.

Comparative Results. Tables 2a and 2b present comparative results between existing FAS methods and the proposed InstructFLIP on the proposed unified evaluation protocol. InstructFLIP consistently outperforms all prior approaches across all evaluation metrics and datasets, demonstrating its strong robustness and adaptability. In terms of HTER, InstructFLIP achieves significant reductions compared to the current SOTA method CFPL [23], with improvements of 37%, 47%, 3.5%, 16%, 28%, and 25% on the M, C, I, O, W, and S datasets, respectively. Beyond HTER, consistent gains in AUC reveal InstructFLIP’s ability to learn robust and discriminative features for distinguishing genuine from spoofed faces. Improvements in TPR@FPR=1% further demonstrate the effectiveness of InstructFLIP in maintaining a low false positive rate while ensuring high detection accuracy, which is essential for practical deployment where both security and user experience are critical.

Table 3: Ablation study for the effectiveness of each component in InstructFLIP, where CB, SB, and Cue denote content branch, style branch, and cue generator, respectively.

CB	SB	Cue	HTER↓	AUC↑	TPR↑@FPR=1%
–	–	–	21.96 (+0.0%)	82.88 (+0.0%)	26.89 (+0.0%)
✓	–	–	14.25 (+35.1%)	90.23 (+8.9%)	32.32 (+20.2%)
–	✓	–	16.25 (+26.0%)	88.27 (+6.5%)	41.92 (+55.9%)
✓	✓	–	13.25 (+39.7%)	91.21 (+10.1%)	60.84 (+126.3%)
✓	✓	✓	12.68 (+42.3%)	93.68 (+13.0%)	65.23 (+142.6%)

Despite the overall strong performance, the results on the CASIA-CeFA (C) dataset remain relatively modest, indicating potential limitations in modeling granular cues. This suggests a direction for future work focused on enhancing the model’s sensitivity to subtle cultural and environmental variations through more expressive and adaptive instruction-tuning strategies. Overall, the results summarized in Table 2c validate that InstructFLIP offers a well-balanced and generalizable solution for robust FAS across diverse domains.

4.3 Ablation Study

In this subsection, we perform a series of ablation studies to assess the effectiveness of the proposed method from multiple perspectives. All ablation experiments are conducted using the meta-domain for training, and the reported results represent the average performance across the seven target datasets, following the unified evaluation protocol described in Section 4.1.

Effectiveness of proposed components. We begin by evaluating the contribution of each component in the proposed InstructFLIP framework through an ablation study, including the content branch (CB), style branch (SB), and cue generation (Cue). The baseline model is constructed by removing all proposed components and relying solely on the given image, following the conventional FAS setup for a fair comparison. The comparative results, summarized in Table 3, illustrate the impact of each module on overall performance.

Introducing the CB to the baseline model resulted in substantial improvements: +35.1% in HTER, +8.9% in AUC, and +20.2% in TPR@FPR=1%. These results demonstrate the critical role of CB, which utilizes fine-grained labels to enhance the model’s ability to effectively capture features that differentiate genuine inputs from spoofing attacks. The contribution of the SB is also evaluated, which further improved performance by +26.0% in HTER, +6.5% in AUC, and +55.9% in TPR@FPR=1%. The considerable boost in TPR@FPR=1% indicates that SB effectively models non-spoofing patterns, thereby reducing overfitting and improving generalization. When both CB and SB were integrated, the model achieved even greater performance gains: +39.7% in HTER, +10.1% in AUC, and +126.3% in TPR@FPR=1%, demonstrating their complementary nature in capturing both spoof-relevant and irrelevant contextual features. Furthermore, the addition of Cue achieved SOTA performance, yielding improvements of +42.3% in HTER, +13.0% in AUC, and an impressive +142.6% increase in TPR@FPR=1%, by leveraging cue supervision to enhance the model’s ability to effectively distinguish samples. These results confirm the effectiveness of each proposed component and the synergistic benefits of their integration in enhancing robustness and generalization.

Table 4: Ablation study to assess the effectiveness of frozen LLMs in InstructFLIP. InstructFLIP[†] replaces LLM in the dual branches with classification heads.

Method	HTER↓	AUC↑	TPR↑@FPR=1%
CFPL [23]	16.15	88.65	53.98
InstructFLIP [†]	14.39	89.97	48.33
InstructFLIP	12.68	93.68	65.23

Table 5: Ablation studies on effectiveness of category diversity. In the table, *Pr.*, *Re.*, M_{2D} , and M_{3D} represent print, replay, 2D mask, and 3D mask, respectively.

Real	<i>Pr.</i>	<i>Re.</i>	M_{2D}	M_{3D}	HTER↓	AUC↑	TPR↑@FPR=1%
1	1	1	1	1	19.11	87.14	39.13
1	2	2	2	1	18.57	87.3	38.3
1	3	3	3	1	12.68	93.68	65.23

Effectiveness of LLMs adaptation. To better understand whether the performance gains observed in Section 4.2 are primarily attributed to the use of frozen large language models (LLMs) or to the incorporation of fine-grained labels from the meta-domain, we conduct an additional experiment with an ablated variant, InstructFLIP[†]. This variant retains the content and style branches but replaces the frozen LLMs with lightweight classification heads. As shown in Table 4, incorporating fine-grained semantic signals already yields performance superior to CFPL, and further integrating LLMs leads to the best overall results. These results point out the complementary roles of structured supervision and language-based reasoning in enhancing model generalization and discriminative capability.

Effectiveness of category diversity. To evaluate the effect of semantic granularity in spoof category annotations, we conduct an ablation study on the category diversity used during training. The results, shown in Table 5, demonstrate the positive impact of incorporating detailed subcategories for print, replay, and 2D mask attacks. As the number of fine-grained classes increases from 1 to 3, the model exhibits consistent improvements across all key evaluation metrics. Specifically, the HTER decreases from 19.11 to 12.68, indicating a substantial reduction in overall error. The AUC improves from 87.14 to 93.68, suggesting enhanced discriminative capability, while the TPR@FPR=1% increases significantly from 39.13 to 65.23, highlighting improved sensitivity under low false positive constraints. These results indicate that richer semantic supervision allows the model to better capture intra-class variations, thereby improving robustness and reliability in real-world use.

Effectiveness of style prompts. In order to investigate the impact of style prompt diversity on model performance, we conduct an ablation study by progressively incorporating different style prompts into InstructFLIP, with the results summarized in Table 6. Starting to use only a single style prompt related to illumination conditions (style1) yields an HTER of 20.65, an AUC of 85.60, and a TPR@FPR=1% of 37.48. Adding an additional prompt reflecting environmental context (style1 + style2) leads to noticeable improvements, reducing the HTER to 19.25, increasing the AUC to 86.71, and raising the TPR@FPR=1% to 45.08. Further incorporating a third prompt associated with camera quality (style1 + style2 + style3)

Table 6: Ablation studies on effectiveness of style prompts.

Style prompt	HTER↓	AUC↑	TPR↑@FPR=1%
Style1	20.65	85.60	37.48
Style1 + Style2	19.25	86.71	45.08
Style1 + Style2 + Style3	12.68	93.68	65.23

Table 7: Question inputs for VLMs.

Type	Question inputs
Content	Which type of spoof is in this image? (1) Real face (2) Photo (3) Poster (4) A4-paper (5) 2D face mask (6) 2D upper-body mask (7) 2D region mask (8) PC screen (9) Pad screen (10) Phone screen (11) 3D mask
Style1	What is the illumination condition in this image? (1) Normal (2) Strong (3) Back (4) Dark
Style2	What is the environment in this image? (1) Indoor (2) Outdoor
Style3	What is the camera quality in this image? (1) Low (2) Medium (3) High

significantly enhances SOTA performance. The results demonstrate the importance of incorporating diverse style prompts to model non-spoof-related variations. By enriching the instruction space, the model becomes more effective at disentangling spoof-relevant and confounding factors, leading to improved generalization and robustness across diverse domain conditions.

4.4 Qualitative Results

In this subsection, we present qualitative results to illustrate the strengths and limitations of the proposed method. The analysis includes successful cases, failure cases, and comparisons with open-source VLMs to provide a clearer understanding of InstructFLIP’s capabilities. To ensure a fair and informative comparison, each input to the VLMs is structured using four prompt components—content, style1, style2, and style3—as defined in Table 7. This design enables the models to develop a comprehensive understanding of each sample’s semantic and contextual attributes.

Illustration of successful cases. Tables 8a and 8b present successful cases processed by the proposed InstructFLIP model. For the true positive sample in Table 8a, the model assigns a fake score of 0.001, indicating a negligible likelihood of misclassifying the live face as a spoof. In contrast, the true-negative sample in Table 8b receives a fake score of 0.998, reflecting high confidence in correctly identifying the spoof. The LLM-generated predictions—including spoof type, illumination, environment, and camera quality—demonstrate InstructFLIP’s ability to generalize to unseen domains by leveraging structured instructional inputs from the meta domain. Furthermore, the cue representations contribute significantly to differentiating between live and spoofed samples, highlighting their importance in the overall decision-making process.

Illustration of failure cases. To better understand the limitations of our model, we analyzed representative failure cases, categorizing them into false positives and false negatives. For the false-positive case illustrated in Table 8c, the model misclassified a spoofed poster

Table 8: Illustration of (a) True-positive, (b) True-negative, (c) False-positive, and (d) False-negative samples predicted by the proposed method. Red indicates incorrect answers.


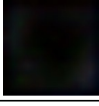


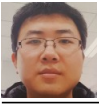
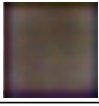
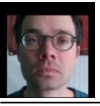
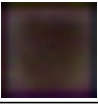


(a) True-positive sample						(b) True-negative sample					
Raw image	Cue map	Fake score	Q-type	Answer	GT	Raw image	Cue map	Fake score	Q-type	Answer	GT
		0.001	Content Style1 Style2 Style3	(1) Real face (3) Back (2) Outdoor (3) High	(1) Real face (3) Back (2) Outdoor (3) High			0.998	Content Style1 Style2 Style3	(8) PC Screen (1) Normal (1) Indoor (2) Medium	(8) PC Screen (1) Normal (1) Indoor (2) Medium
(c) False-positive sample						(d) False-negative sample					
Raw image	Cue map	Fake score	Q-type	Answer	GT	Raw image	Cue map	Fake score	Q-type	Answer	GT
		0.055	Content Style1 Style2 Style3	(1) Real face (3) Back (1) Indoor (2) Medium	(3) Poster (1) Normal (1) Indoor (2) Medium			0.851	Content Style1 Style2 Style3	(8) PC screen (3) Back (1) Indoor (3) High	(1) Real face (3) Back (1) Indoor (2) Medium

Table 9: Comparison of model responses. Red indicates incorrect answers and gray represents indirect or ambiguous responses.

Raw image	Question type	Response			GT
		InstructBLIP	ChatGPT-4o	InstructFLIP	
	Content	(1) Real face	I cannot determine spoof types from an uploaded image. However, you can analyze this image using algorithms trained for Face Anti-Spoofing tasks. Let me know if you need guidance on setting up such an analysis.	(9) Pad screen	(9) Pad screen
	Illumination	(1) Normal	I cannot visually analyze the illumination condition from the image. However, if you have an algorithm to evaluate the illumination or would like to process the image computationally, I can guide you through that.	(1) Normal	(1) Normal
	Content	(1) Real face	I cannot determine spoof types from an uploaded image. However, you can analyze this image using algorithms trained for Face Anti-Spoofing tasks. Let me know if you need guidance on setting up such an analysis.	(1) Real face	(1) Real face
	Environment	(1) Indoor	The image appears to be taken in an outdoor environment, as suggested by the presence of natural sunlight, reflections in the sunglasses, and background elements such as trees or open space.	(2) Outdoor	(2) Outdoor

as a real face, likely due to difficulty in distinguishing subtle material cues like gloss and texture. It also misjudged backlighting as Normal, indicating sensitivity to slight luminance variations, though it accurately predicted other style attributes such as environment and camera quality. For the false-negative case shown in Table 8d, a real face was classified as a spoofed PC screen, possibly due to overfitting to reflective surface patterns. While illumination was correctly predicted, the camera quality was overestimated, revealing inconsistency in assessing image sharpness.

Comparison with open VLMs. We qualitatively compare our InstructFLIP with InstructBLIP and GPT-4o using both content and style-based instructions on a fake and a live images, respectively, as shown in Table 9. For content-based instructions, InstructFLIP consistently outperforms the other VLMs, accurately identifying fake samples as Pad screens and correctly recognizing live faces. In contrast, InstructBLIP frequently misclassifies spoofed faces as real, indicating a limited capacity to capture spoofing cues. GPT-4o, although capable of generating general explanations, refrains from making explicit predictions and instead suggests computational techniques, reducing its utility in this task. For style-based instructions, InstructFLIP provides strong performance across diverse conditions, accurately predicting illumination, camera quality, and environment. While InstructBLIP performs reasonably on simpler attributes like illumination, it struggles with more nuanced aspects

such as environment, occasionally misclassifying Outdoor as Indoor. GPT-4o avoids direct responses in the style setting, revealing a lack of task-specific grounding. These findings show InstructFLIP's adaptability across both instruction types and underscore the importance of efficient contextual understanding in VLMs.

5 Conclusion

We present InstructFLIP, a novel instruction approach for FAS, featuring a unified architecture that decouples content and style representations. By leveraging a dual-branch design, InstructFLIP effectively captures spoof-related patterns from a meta domain, ensuring robust generalization across diverse domains. We employ the Q-Former to encode semantic information from content and style prompts, allowing the model to comprehensively understand diverse spoofing types and environmental conditions. We also propose a generalized query fusion strategy and additional cue maps that enhance FAS performance by efficiently capturing domain-invariant characteristics. Experimental results validate InstructFLIP's efficacy in achieving strong generalization, marking a significant step toward practical, adaptable FAS solutions in real-world applications. Future work may explore extending this framework to other visual tasks where robustness across domains remains challenging, further advancing instruction-driven generalization.

6 Acknowledgments

This work is partially supported by the National Science and Technology Council, Taiwan, under Grant: NSTC-112-2628-E-002-033-MY4, and was financially supported in part by the Center of Data Intelligence: Technologies, Applications, and Systems, National Taiwan University (Grants: 114L900901/114L900902/114L900903), from the Featured Areas Research Center Program within the framework of the Higher Education Sprout Project by the Ministry of Education, Taiwan.

References

- [1] Xiaoming Liu, Amin Jourabloo*, Yaojie Liu*. 2018. Face De-Spoofing: Anti-Spoofing via Noise Modeling. In *ECCV*. 297–315.
- [2] Zinelabinde Boulkenafet, Jukka Komulainen, Lei Li, Xiaoyi Feng, and Abdenour Hadid. 2017. OULU-NPU: A mobile face presentation attack database with real-world variations. In *IEEE FG*. IEEE, 612–618.
- [3] Shunxin Chen, Ajian Liu, Junze Zheng, Jun Wan, Kailai Peng, Sergio Escalera, and Zhen Lei. 2025. Mixture-of-Attack-Experts with Class Regularization for Unified Physical-Digital Face Attack Detection. In *AAAI*. 2195–2203.
- [4] Ivana Chingovska, André Anjos, and Sébastien Marcel. 2012. On the effectiveness of local binary patterns in face anti-spoofing. In *IEEE BIOSIG*. IEEE, 1–7.
- [5] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. *JMLR* 25, 70 (2024), 1–53.
- [6] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning. In *NeurIPS*.
- [7] Xin Dong, Hao Liu, Weiwei Cai, Pengyuan Lv, and Zekuan Yu. 2021. Open Set Face Anti-Spoofing in Unseen Attacks. In *ACMMM*. Association for Computing Machinery, 4082–4090.
- [8] Zhekai Du, Jingjing Li, Lin Zuo, Lei Zhu, and Ke Lu. 2022. Energy-Based Domain Generalization for Face Anti-Spoofing. In *ACMMM*. Association for Computing Machinery, 1749–1757.
- [9] Xinxu Ge, Xin Liu, Zitong Yu, Jingang Shi, Chun Qi, Jie Li, and Heikki Kälviäinen. 2024. Diffias: face anti-spoofing via generative diffusion models. In *ECCV*. Springer, 144–161.
- [10] Anjith George, Zohreh Mostaani, David Geissenbuhler, Olegs Nikisins, André Anjos, and Sébastien Marcel. 2019. Biometric face presentation attack detection with multi-channel convolutional neural network. *IEEE Trans. Inf. Forensics Secur.* 15 (2019), 42–55.
- [11] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. 2022. Open-vocabulary Object Detection via Vision and Language Knowledge Distillation. In *ICLR*.
- [12] Jiabao Guo, Huan Liu, Yizhi Luo, Xueli Hu, Hang Zou, Yuan Zhang, Hui Liu, and Bo Zhao. 2024. Style-conditional Prompt Token Learning for Generalizable Face Anti-spoofing. In *ACMMM*. Association for Computing Machinery, 994–1003.
- [13] Xueli Hu, Huan Liu, Haocheng Yuan, Zhiyang Fu, Yizhi Luo, Ning Zhang, Hang Zou, Jianwen Gan, and Yuan Zhang. 2024. Fine-Grained Prompt Learning for Face Anti-Spoofing. In *ACMMM*. 7619–7628.
- [14] Hsin-Ping Huang, Deqing Sun, Yaojie Liu, Wen-Sheng Chu, Taihong Xiao, Jinwei Yuan, Hartwig Adam, and Ming-Hsuan Yang. 2022. Adaptive transformers for robust few-shot cross-domain face anti-spoofing. In *ECCV*. Springer, 37–54.
- [15] Pei-Kai Huang, Cheng-Hsuan Chiang, Tzu-Hsien Chen, Jun-Xiong Chong, Tyng-Luh Liu, and Chiou-Ting Hsu. 2024. One-Class Face Anti-spoofing via Spoof Cue Map-Guided Feature Learning. In *CVPR*. 277–286.
- [16] Xun Huang and Serge Belongie. 2017. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*. 1501–1510.
- [17] Yunpei Jia, Jie Zhang, Shiguang Shan, and Xilin Chen. 2020. Single-side domain generalization for face anti-spoofing. In *CVPR*. 8484–8493.
- [18] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*. PMLR, JMLR.org, 19730–19742.
- [19] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*. PMLR, 12888–12900.
- [20] Ajian Liu, Xuan Li, Jun Wan, Yanyan Liang, Sergio Escalera, Hugo Jair Escalante, Meysam Madadi, Yi Jin, Zhuoyuan Wu, Xiaogang Yu, et al. 2021. Cross-ethnicity face anti-spoofing recognition challenge: A review. *IET Biom.* 10, 1 (2021), 24–43.
- [21] Ajian Liu, Hui Ma, Junze Zheng, Haocheng Yuan, Xiaoyuan Yu, Yanyan Liang, Sergio Escalera, Jun Wan, and Zhen Lei. 2024. FM-CLIP: Flexible Modal CLIP for Face Anti-Spoofing. In *ACMMM*. Association for Computing Machinery, 8228–8237.
- [22] Ajian Liu, Zichang Tan, Jun Wan, Sergio Escalera, Guodong Guo, and Stan Z Li. 2021. Casia-surf cefa: A benchmark for multi-modal cross-ethnicity face anti-spoofing. In *WACV*. 1179–1187.
- [23] Ajian Liu, Shuai Xue, Jianwen Gan, Jun Wan, Yanyan Liang, Jiankang Deng, Sergio Escalera, and Zhen Lei. 2024. CFPL-FAS: Class Free Prompt Learning for Generalizable Face Anti-spoofing. In *CVPR*. 222–232.
- [24] Si-Qi Liu, Qirui Wang, and Pong C Yuen. 2024. Bottom-up domain prompt tuning for generalized face anti-spoofing. In *ECCV*. Springer, 170–187.
- [25] Yaojie Liu, Amin Jourabloo, and Xiaoming Liu. 2018. Learning Deep Models for Face Anti-Spoofing: Binary or Auxiliary Supervision. In *CVPR*. 389–398.
- [26] Ilya Loshchilov and Frank Hutter. 2017. Decoupled Weight Decay Regularization. In *ICLR*.
- [27] Keyurkumar Patel, Hu Han, and Anil K Jain. 2016. Secure face unlock: Spoof detection on smartphones. *IEEE Trans. Inf. Forensics Secur.* 11, 10 (2016), 2268–2283.
- [28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*. PMLR, 8748–8763.
- [29] Leslie N Smith and Nicholay Topin. 2019. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial intelligence and machine learning for multi-domain operations applications*. SPIE.
- [30] Koushik Srivatsan, Muzammal Naseer, and Karthik Nandakumar. 2023. Flip: Cross-domain face anti-spoofing with language guidance. In *ICCV*. 19685–19696.
- [31] Joel Stehouwer, Amin Jourabloo, Yaojie Liu, and Xiaoming Liu. 2020. Noise modeling, synthesis and classification for generic object anti-spoofing. In *CVPR*. 7294–7303.
- [32] Yiyou Sun, Yaojie Liu, Xiaoming Liu, Yixuan Li, and Wen-Sheng Chu. 2023. Re-thinking domain generalization for face anti-spoofing: Separability and alignment. In *CVPR*. 24563–24574.
- [33] Chien-Yi Wang, Yu-Ding Lu, Shang-Ta Yang, and Shang-Hong Lai. 2022. Patchnet: A simple face anti-spoofing framework via fine-grained patch recognition. In *CVPR*. 20281–20290.
- [34] Guoqing Wang, Hu Han, Shiguang Shan, and Xilin Chen. 2020. Cross-domain face presentation attack detection via multi-domain disentangled representation learning. In *CVPR*. 6678–6687.
- [35] Guoqing Wang, Hu Han, Shiguang Shan, and Xilin Chen. 2020. Unsupervised adversarial domain adaptation for cross-domain face presentation attack detection. *IEEE Trans. Inf. Forensics Secur.* 16 (2020), 56–69.
- [36] Yu-Chun Wang, Chien-Yi Wang, and Shang-Hong Lai. 2022. Disentangled representation with dual-stage feature learning for face anti-spoofing. In *WACV*. 1955–1964.
- [37] Zhuo Wang, Zezheng Wang, Zitong Yu, Weihong Deng, Jiahong Li, Tingting Gao, and Zhongyuan Wang. 2022. Domain generalization via shuffled style assembly for face anti-spoofing. In *CVPR*. 4123–4133.
- [38] Di Wen, Hu Han, and Anil K Jain. 2015. Face spoof detection with image distortion analysis. *IEEE Trans. Inf. Forensics Secur.* 10, 4 (2015), 746–761.
- [39] Jianwei Yang, Zhen Lei, and Stan Z Li. 2014. Learn convolutional neural network for face anti-spoofing. *arXiv preprint arXiv:1408.5601* (2014).
- [40] Xiao Yang, Wenhan Luo, Linchao Bao, Yuan Gao, Dihong Gong, Shibao Zheng, Zhifeng Li, and Wei Liu. 2019. Face Anti-Spoofing: Model Matters, so Does Data. In *CVPR*. 3502–3511.
- [41] Zitong Yu, Yunxiao Qin, Xiaobai Li, Chenxu Zhao, Zhen Lei, and Guoying Zhao. 2022. Deep learning for face anti-spoofing: A survey. *IEEE TPAMI* 45, 5 (2022), 5609–5631.
- [42] Zitong Yu, Chenxu Zhao, Zezheng Wang, Yunxiao Qin, Zhuo Su, Xiaobai Li, Feng Zhou, and Guoying Zhao. 2020. Searching Central Difference Convolutional Networks for Face Anti-Spoofing. In *CVPR*. 5295–5305.
- [43] Haixiao Yue, Keyao Wang, Guosheng Zhang, Haocheng Feng, Junyu Han, Errui Ding, and Jingdong Wang. 2023. Cyclically disentangled feature translation for face anti-spoofing. In *AAAI*, Vol. 37. 3358–3366.
- [44] Guosheng Zhang, Keyao Wang, Haixiao Yue, Ajian Liu, Gang Zhang, Kun Yao, Errui Ding, and Jingdong Wang. 2025. Interpretable Face Anti-Spoofing: Enhancing Generalization with Multimodal Large Language Models. In *AAAI*. 9896–9904.
- [45] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. 2016. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Process. Lett.* 23, 10 (2016), 1499–1503.
- [46] Renrui Zhang, Wei Zhang, Rongyao Fang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. 2022. Tip-adapter: Training-free adaption of clip for few-shot classification. In *ECCV*. Springer, 493–510.
- [47] Shifeng Zhang, Ajian Liu, Jun Wan, Yanyan Liang, Guodong Guo, Sergio Escalera, Hugo Jair Escalante, and Stan Z Li. 2020. Casia-surf: A large-scale multi-modal benchmark for face anti-spoofing. *IEEE Trans. Biom. Behav. Identity Sci.* 2, 2 (2020), 182–193.
- [48] Shifeng Zhang, Xiaobo Wang, Ajian Liu, Chenxu Zhao, Jun Wan, Sergio Escalera, Hailin Shi, Zezheng Wang, and Stan Z Li. 2019. A dataset and benchmark for large-scale multi-modal face anti-spoofing. In *CVPR*. 919–928.
- [49] Yuanhan Zhang, Zhenfei Yin, Yidong Li, Guojun Yin, Junjie Yan, Jing Shao, and Ziwei Liu. 2020. CelebA-Spoof: Large-Scale Face Anti-Spoofing Dataset with Rich Annotations. In *ECCV*. 70–85.

- [50] Zhiwei Zhang, Junjie Yan, Sifei Liu, Zhen Lei, Dong Yi, and Stan Z Li. 2012. A face antispoofing database with diverse attacks. In *IEEE ICB*. IEEE, 26–31.
- [51] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Conditional prompt learning for vision-language models. In *CVPR*. 16816–16825.
- [52] Qianyu Zhou, Ke-Yue Zhang, Taiping Yao, Xuequan Lu, Ran Yi, Shouhong Ding, and Lizhuang Ma. 2023. Instance-Aware Domain Generalization for Face Anti-Spoofing. In *CVPR*. 20453–20463.
- [53] Qianyu Zhou, Ke-Yue Zhang, Taiping Yao, Ran Yi, Shouhong Ding, and Lizhuang Ma. 2022. Adaptive Mixture of Experts Learning for Generalizable Face Anti-Spoofing. In *ACMMM*. Association for Computing Machinery, 6009–6018.
- [54] Qianyu Zhou, Ke-Yue Zhang, Taiping Yao, Ran Yi, Kekai Sheng, Shouhong Ding, and Lizhuang Ma. 2022. Generative domain adaptation for face anti-spoofing. In *ECCV*. Springer, 335–356.
- [55] Hang Zou, Chenxi Du, Hui Zhang, Yuan Zhang, Ajian Liu, Jun Wan, and Zhen Lei. 2024. La-SoftMoE CLIP for Unified Physical-Digital Face Attack Detection. In *IJCB*. 1–11. doi:10.1109/IJCB62174.2024.10744523