

GM-DF: Generalized Multi-Scenario Deepfake Detection

Yingxin Lai, Zitong Yu, Jing Yang, Bin Li, Xiangui Kang, Linlin Shen

Abstract—Existing face forgery detection usually follows the paradigm of training models in a single domain, which leads to limited generalization capacity when unseen scenarios and unknown attacks occur. In this paper, we elaborately investigate the generalization capacity of deepfake detection models when jointly trained on multiple face forgery detection datasets. We first find a rapid degradation of detection accuracy when models are directly trained on combined datasets due to the discrepancy across collection scenarios and generation methods. To address the above issue, a Generalized Multi-Scenario Deepfake Detection framework (GM-DF) is proposed to serve multiple real-world scenarios by a unified model. First, we propose a hybrid expert modeling approach for domain-specific real/forgery feature extraction. Besides, as for the commonality representation, we use CLIP to extract the common features for better aligning visual and textual features across domains. Meanwhile, we introduce a masked image reconstruction mechanism to force models to capture rich forged details. Finally, we supervise the models via a domain-aware meta-learning strategy to further enhance their generalization capacities. Specifically, we design a novel domain alignment loss to strongly align the distributions of the meta-test domains and meta-train domains. Thus, the updated models are able to represent both specific and common real/forgery features across multiple datasets. In consideration of the lack of study of multi-dataset training, we establish a new benchmark leveraging multi-source data to fairly evaluate the models' generalization capacity on unseen scenarios. Both qualitative and quantitative experiments on five datasets conducted on traditional protocols as well as the proposed benchmark demonstrate the effectiveness of our approach. The codes will be available on <https://github.com/laiyingxin2/GM-DF>.

Index Terms—face forgery detection, domain generalization, meta-learning, CLIP, masked image reconstruction.

1 INTRODUCTION

ADVANCEMENTS in deep learning have facilitated the creation of face forgery mechanisms [1], [2], [3], [4], [5], [6]. These techniques simplify the generation of highly realistic forged face images, posing risks to both political and personal reputations and giving rise to significant social challenges. Consequently, the development of detection methods to mitigate these risks is imperative. To alleviate discrepancies among various face forgery detection datasets, some researchers have adopted a specific approach. They treat the task of detecting forged faces as a binary classification task, utilizing existing deep convolutional neural

networks to categorize the data into two distinct classes: real and forged. The primary goal of these investigations is to identify and extract common features to address the challenge of feature discrepancies. Several approaches have been proposed to tackle this issue, including the use of noise as a form of supervision [7], [8], the incorporation of frequency domain information [9], [9], [10], and the application of reconstruction techniques to gain insights into the distribution of authentic samples [11], [12].

However, despite the remarkable accuracy and precision attained by these models when applied in a cross-domain setting, their effectiveness remains heavily dependent upon the training process conducted on only one dataset. The initial strategy involves training a baseline model on the combined datasets. However, the results shown in Figure 1 indicate that direct training within combined datasets easily leads to generalization drops. The main reasons behind this might be the variances in forgery techniques, capturing circumstances, forgery methods, and hardware across various domains. In light of the continuous emergence of manipulated facial datasets, it is imperative to integrate and simultaneously train using different accessible data sources.

But if two face forgery detection datasets with different distributions are directly merged and used for training, the problem of domain conflict will inevitably be encountered. For example, as shown in Figure 1, the merging of Celeb-DF(V2) [15] which the original FF++ [14] dataset, suffers from degradation of accuracy from 75.49% to 67.19%. Therefore, the previous paradigm of single dataset training and testing does not work well on multiple domains, and the direct merging of individual datasets does not improve the generalization ability of the model well. With the increasing

Manuscript received May 2024. Corresponding author: Zitong Yu (email: zitong.yu@ieee.org).

This work was supported by Open Fund of National Engineering Laboratory for Big Data System Computing Technology (Grant No. SZU-BDSC-OF2024-02) and National Natural Science Foundation of China under Grant 62306061.

Y. Lai and J. Yang are with the School of Computing and Information Technology, Great Bay University, Dongguan 523000, China.

Z. Yu is with the School of Computing and Information Technology, Great Bay University, Dongguan 523000, China, and National Engineering Laboratory for Big Data System Computing Technology, Shenzhen University, Shenzhen 518060, China

B. Li is with the Guangdong Key Laboratory of Intelligent Information Processing, Shenzhen Key Laboratory of Media Security, Guangdong Laboratory of Artificial Intelligence and Digital Economy (SZ), Shenzhen Institute of Artificial Intelligence and Robotics for Society, Shenzhen University, Shenzhen 518060, China.

X. Kang is with the Guangdong Key Laboratory of Information Security, and the School of Computer Science and Engineering, Sun Yat-sen University, Guangzhou 510080, China

L. Shen is with Computer Vision Institute, School of Computer Science & Software Engineering, Shenzhen Institute of Artificial Intelligence and Robotics for Society, Guangdong Key Laboratory of Intelligent Information Processing, and National Engineering Laboratory for Big Data System Computing Technology, Shenzhen University, Shenzhen 518060, China.

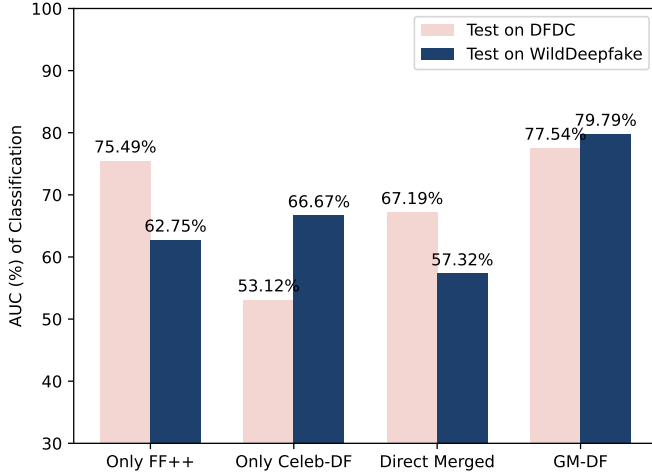


Fig. 1: Challenges in training a detector from multiple datasets. The generalization capacity of the baseline Xception [13] trained on FF++ [14]&Celeb [15] datasets drops sharply while the proposed method GM-DF benefits obviously from multi-dataset training.

number of various face forgery detection datasets, *how to effectively train a unified detector on multiple widely differentiated datasets is worth exploring. The solutions behind the problem might benefit the development of forgery foundation models.*

Based on the above observations, in this paper, we propose a unified face forgery detection framework to solve the multi-dataset conflict problem, and our model is orthogonal to existing methods. We discover a novel insight: data conflict might be caused by ignoring the domain-specific features of the datasets. In order to enhance the models' generalization ability with the increasing number of datasets, we design a hybrid expert modeling approach to extract the domain-specific features while leveraging image-text alignment and masked image reconstruction mechanism to extract common real/forgery features across domains. Finally, we supervise the models via a domain-aware meta-learning strategy. we design the novel domain alignment loss to strongly align the distributions of the meta-test domains and meta-train domains. Thus, the updated models are able to represent both specific and common real/forgery features across multiple datasets.

Extensive experiments are conducted on five public autopilot datasets, including FF++ [14] Celeb-DF(V2) [15], WildDeepfake [16], DFDC [17] and the fake face dataset generated by diffusion DFF [18] to study the problem of data conflict in each domain or merged domains. Towards the era of large-scale multi-dataset training and testing, we establish a novel benchmark with five mainstream datasets, and the results show that the proposed models have strong generalization ability. Our contributions are as follows:

- We are the first to comprehensively investigate the multi-dataset training task for face forgery detection, and establish a new benchmark leveraging multi-dataset data to fairly evaluate the models' generalization capacity.
- We propose a hybrid expert modeling approach for domain-specific real/forgery feature extraction. We

also propose to represent common features via simultaneously aligning visual and textual features, and reconstructing masked faces across domains.

- We supervise the models via a domain-aware meta-learning strategy with a novel domain alignment loss.
- The proposed method achieve state-of-the-art performance on both traditional protocols as well as the proposed benchmark.

2 RELATED WORK

In this section, we briefly describe deepfake detection, vision language models, and joint training on multiple datasets.

2.1 Face Forgery Detection

Recently, face forgery detection has received extensive attention from researchers due to the great threat to security and privacy. Previous methods [13], [19], [20], [21] treat face forgery detection as a binary classification problem in intra-dataset testing and has achieved very satisfactory performance. However, only relying on binary classification supervision, deepfake detectors easily overfit the training data. Thus, people began to turn to cross-domain performance exploration, F³-Net [10] combines with the frequency domain information to extract the subtle differences between real and fake pictures, proved the effectiveness of the frequency domain in forgery detection artifact recognition. Similarly, SPSL [22] proposes a frequency-based phase spectral analysis method. Face X-ray [23] detects generated images by picture mixing boundaries, and DADF [24] leverages vision foundation model for robust forgery localization. PCL [25] improves the supervisory performance by learning the inconsistency between the forged/ neighboring regions and learning the commonality from real samples while reconstructing real samples. SLADD [26] combines data augmentation and face blending to improve the generalization ability. M-FAS [27] and [5] established a unified face forgery detection system.

Although these methods substantially improve the generalization ability, they are still limited by the common features and the specific forgery patterns in the training set, which aggravates the data Conflicts.

2.2 Vision Language Models

Visual-language models are rich in multimodal feature representations and show surprising generalization performance in downstream tasks. [28] proposes an adaptive approach to CLIP [29] modules without training that performs state-of-the-art small-sample classification tasks on ImageNet. CoOp [30] aims to introduce a learnable Prompt approach to better adapt powerful and generalized a priori of visual-linguistic models to downstream tasks. OpenCLIP [31] The model integrates the cross-modal capabilities of text encoder with the generative abilities of the pre-trained language model BART, resulting in a strengthened text encoder for language backbone. Lit [32] utilizes multimodal pre-trained models to improve graphic alignment. Flamingo [33] predicts the next text token based on the previous text and the visual Token, thus better introducing visual information for text creation. LLAVA [34] proposes a command

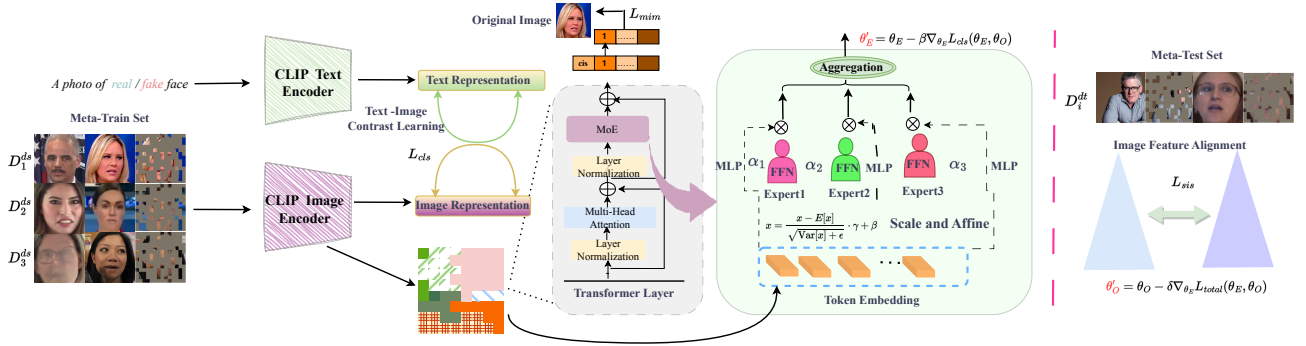


Fig. 2: The framework of the proposed method. It integrates meta-learning modeling with image-text contrastive learning. It comprises three pivotal components: Dataset-Embedding Generator (DEG) and a Multi-Dataset Representation (MDP), as well as a Meta-Domain-Embedding Optimizer(MDEO). Firstly, the DEG incorporates a Dataset Information Layer (DIL) and a dynamic text feature affine aimed at mapping discriminative features unique to each domain, and the second part MDP is the face mask image modeling (MIM) reconstruction module, which provides additional detail information for the global features of CLIP. To consider the difference between each domain, we propose to use the higher-order statistical features in Domain Alignment (DA) loss to constrain the feature distribution. In this process, MDEO was used to optimize the learned two features.

optimization technique for vision. BLIP2 [35] designs Q-Former to bridge between visual and linguistic models by connecting temporal and linguistic features. Although these methods achieve good generalization performance in downstream tasks, they face the challenges of high computational effort and complexity. In addition, they are mostly applied to face forgery detection, where lack of robustness remains a problem.

2.3 Joint Training on Multiple Datasets

For traditional image tasks such as target detection [36], [37] and semantic segmentation [38], [39], due to the different dataset class labels and fine-grained cross-dataset difference, they result in poor generalization when directly fused data for training. Some researchers have begun to study the data federation [40], [41], [42], [43], [44]. Dai et al. [40] combine multiple self-attention mechanisms sequentially to unify the target detection head [44] and relabeling the instances of disaggregation to perform the alignment operation on to the images significantly improve the generalization ability of the model. Wang et al. [42] trains a generalized object detector by incorporating different supervised signals, eliminating the need to model differences across data. Zhao et al. [45] propose a pseudo-labeling method that is tuned for specific situations, showing that a unified detector trained on multiple datasets can outperform each detector trained on a specific dataset. Although recent works on the generic image classification task using multi-domain data for training have been partially investigated, it has not been explored in the field of face forgery detection. Moreover, different training domains are not equally important due to variant environments, media quality, and attack types. Such biased and imbalanced data from different domains makes this task challenging.

3 METHODOLOGY

The framework of the proposed GM-DF is shown in Figure 2, which contains a Dataset-Embedding Generator (DEG)

and a Multi-dataset Representation (MDP), as well as a Meta-Domain-Embedding Optimizer (MDEO). The DEG pay attention to information that is unique to the dataset, the MDP focuses on learning more fine-grained, local relational features of forged patterns, whereas the MDEO achieves its functionality by modeling the relationships between universal information and dataset embedding. For better understanding, we provide some brief details before outlining the framework architecture.

3.1 Preliminary

To solve the problem of poor cross-domain performance for multiple scenarios and datasets, we first assume that this is due to domain differences. As shown in Table 1, different datasets have different collection scenarios and forgery methods. Currently, there is a trend towards diversification in sources of data for facial forgery detection. Figure 5 highlights distinct differences among various datasets. For instance, the DFDC [17] dataset exhibits a prevalence of green backgrounds, WDF [16] images convey an impression of magnification, DFF [18] tends towards an artistic mode of photography, while the forgeries in Celeb-DF appear relatively homogeneous, and FF++ [14] features facial representations with rich attributes. Mixing these datasets may lead to model learning biases due to their inherent disparities.

We also observe that current deepfake detectors [10], [12], [13], [46] usually focus on representing common patterns. As shown in Figure 3, the distribution of feature differences in each dataset is relatively small, indicating that the model learned some common features but ignored the specific features of each domain. The utilization of frequency domain information for detection is a widely employed technique. As depicted in Figure 4, these methods merely capture singular counterfeit and learning patterns. The consistent frequency domain visualizations across various datasets underscore the imperative nature of learning dataset characteristics.

It is also worth noting that the DFF [18] dataset generated by diffusion method is also not very different from

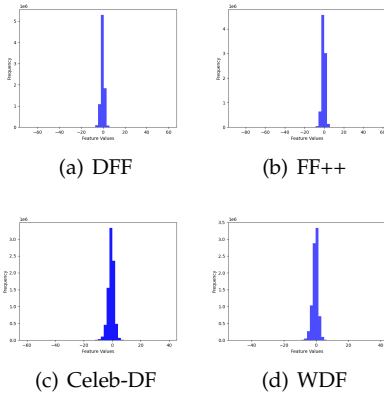


Fig. 3: Histograms of feature values in a randomly selected channel, where features are computed from the block of a convolution based on Xception [13] trained on the dataset of four domains [14], [15], [16], [18].

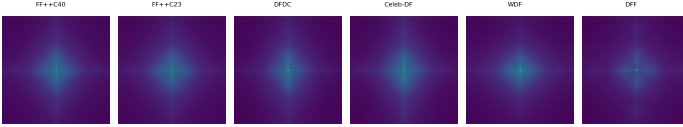


Fig. 4: The commonly used frequency domain detection model M2TR’s [48] frequency domain visualization on the FF++ c40 [14], FF++ c23 [14], DFDC [17], Celeb-DF [15], WildDeepfake [16], and DFF [18] datasets.

the other datasets, leading to large differences in their domains and thus the cause of domain conflicts, so we would like to set up a more specific model that reduces domain conflicts, learns more characteristic forgery features after mapping to the feature space, and at the same time can be well generalized to catch the differences between real and fake images, since forgery patterns are usually hidden in low-level details. Therefore, we refer to the principle of Adaptive Risk Minimization [47], which aims at co-optimal solutions in multiple domains. Specifically, here, we describe our adaptive modelling. Divided into N source domains $D = \{d_{s1}, d_{s2} \dots d_{sn}\}$ represent source face forgery datasets and M target domains. Define $D_t = \{d_{t1}, d_{t2}\}$ where each domain has input and label. Using $x \in X$ and $y \in Y$ as input and label, we may define the source domains as $D_s = \{x_{s^i}, y_{s^i}\}$. To simulate real-world cross-domain challenges by mimicking test-time adaptation (i.e., adjusting prior to prediction), we use characteristic domain weights in the inner loop to learn information unique to each domain, and reconstruction learning and distributional approximation in the outer loop to allow the model to learn the differences between real and fake images.

3.2 Dataset–Embedding Generator

After training use the underlying source domain dataset, it is usually possible to extract a large number of visual features that match the characteristics of the domain. However, when confronted with unseen scenarios and unknown forgery categories, models (e.g., Xception) usually have poor feature generalization (see Figure 1). This is mainly due

to the significant semantic differences between the forgery patterns of the new category and those in the underlying dataset.

For example, when a model processes an image of a face collected under a curtain, it may incorrectly misinterpret features such as the eyes and nose of the face as features of the forgery image under the curtain. This is because there may be some false pictures under the curtain in the underlying dataset, leading to confusion in the model’s learning process. This situation prevents the model from correctly recognizing the forgery images in certain environments or situations. To mitigate this problem, we explore additional semantic information cues to guide the visual feature network to obtain rich and flexible semantic features.

Specifically we use Vit as the foundation model for fine-tuning due to its unbiasedness for each category of both real and fake images and language modeling’s potential. This module follows the Mixture of Experts (MoE) [49] network structure to build a mixture of expert layers to learn domain-invariant features; unlike the domain-specific module we propose based on this, we use N independent experts. Each residual block consists of a Dataset Information Layer and an Multilayer Perceptron (MLP). Since the domain-specific embedding is much smaller than the normal backbone, it can be used if there is a low additional computational cost and restrain the trends of the overfit, experts from various domains carried out the process to extract domain-invariant and domain-specific features as follows:

$$F(x) = F_{\theta}(x) + \Delta F_{\theta}^n(x) \quad (1)$$

Here, $F_{\theta}(x)$ represents the original function that is shared by all source domains to learn the common domain-invariant features. ΔF_{θ}^n adaptively extracts the discriminative and unique domain-specific features.

Although existing works show that the activations in different transformer blocks contribute to the stability of the training, the diversity of the individual domains is sacrificed in the case of multi-domain training. So we model each expert in MoE layer via introducing a new Dataset Information Layer (DIL) with domain-specific parameters. Unlike the fixed gain and bias in LayerNorm, we add skip connections and then scale the function by a learnable parameter called the domain weights and initialize it to 0 at the outset. The signal propagates as follows:

$$x_{i+1} = \alpha_i * \text{Sublayer}(x_i), \quad (2)$$

where $\text{Sublayer} \subseteq \{\text{self-attention, feed-forward}\}$, Then, we compute the gain and bias with respect to the learned prompt vector w . where α_i is the learned residual domain weight parameter. At the initial time, all the α_i are initial to zero; the network represents a constant function, at which point dynamic equidistance is directly satisfied. Then the model gradually learns the specific features corresponding to each domain. In order to allow models to learn their respective domain-specific knowledge through the parameters and to dynamically generate them in real time according to different instances, we use the learned prompt vector to perform an affine to the normalized input features based on VPT [50]. More precisely, given domain d_i^t and prompt

feature vector $p = [v_1, v_2, v_3 \dots, v_M]$ where M is the dimensionality of the learnable prompt vector, we derive a MLP layer $h(\cdot)$ to the specific feature

$$\text{DIL}(v, p) = h(p) \cdot x_{i+1} \quad (3)$$

The entire transformer block can be formalized as follows:

$$\begin{aligned} x_0 &= \text{LayerNormal}_{\text{att}_i}(x), x = \text{MHA}(x_0) + x, \\ x'_0 &= \text{DIL}_{\text{moe}_i}(x), x = \text{MoE}(x'_0) + x. \end{aligned} \quad (4)$$

The input features first go through the original transformer layernorm $\text{LayerNormal}_{\text{att}_i}$ as well as the multi-head attention MHA, and then through the various expert modules $\text{DIL}_{\text{moe}_i}$ and MoE.

3.3 Multi-Dataset Representation

After obtaining the domain embedding and expert views, we calculate the scaled dot-product attention and mark it as the expert views, which is formulated as [11]:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V, \quad (5)$$

where Q denotes the query, K denotes the key, \sum denotes the value of the input embedding, and the scale factor of d_k is the key of dimension. Here we compute the attention score containing the task information Setting Q and K as the

$$Q = K = \text{Concat}(\Delta F_{\theta_1}(x), \Delta F_{\theta_2}(x), \dots, \Delta F_{\theta_N}(x)) \in \mathbb{R}^{1 \times N}, \quad (6)$$

where Concat denotes the operation that stacks vectors into a matrix. V is a matrix stacked by expert views. We make the summation of the expert views to obtain the task-specific aggregated expert view.

Image-text pairs can learn semantic feature representations of face forgeries about specifics, but they may not be able to capture the details. Inspired by previous study on forgery face reconstruction properties [12], [46], [51] and to improve face detail representation, we add a mask image modeling (MIM) [52] task that masks a number of patches of the input image and predicts their visual tokens. Commonly used, typical low-level visual tasks mask the image to capture low-level details and offer semantic information. With the learned representations, the reconstruction difference of real and fake faces significantly differs in distribution.

Given an input image X , we begin by dividing it into N patches denoted as $\{x_1, x_2, x_3, \dots, x_n\}$, where n represents the total number of patches. Subsequently, we adopt a stochastic masking approach, referred to as [53] to apply masks to a subset of M patches. This process results in a modified image X' , expressed as $X' = \{x_1, x'_{m_2}, x'_{m_3}, \dots, x_n\}$. Here, x'_{m_2} means that the second one is replaced by a mask. Next, we feed the masked images into a shared Transformer architecture, yielding a set of hidden vectors $\{h'_{\text{cls}}, h'_1, h'_2, \dots, h'_N\}$. Leveraging the knowledge encapsulated in these hidden vectors, we proceed to predict the masked regions $\{x'_{m_i} \mid m_i \in M\}$ and simultaneously perform direct pixel-level predictions.

To optimize memory consumption, a Gumbel-Softmax Variational Autoencoder [54] is employed. Each image block

Algorithm 1 Training for Meta Deepfake Detection

Input D : data of multi-source domains; δ, β : learning rates;
Initialize: θ_E, θ_O
while not converged **do**
 Sample $N - 1$ domains as meta-train set D_i^{ds} and the remaining domain as meta-test set D_i^{dt} .
 for each D_i^{ds} **do**
 Evaluate loss L_{cls} on D_i^{ds}
 Update θ_E by: $\theta'_E \leftarrow \theta_E - \beta \frac{\sigma L_{cls}(f(\theta_E, \theta_O))}{\sigma \theta_E}$
 end for
 Update θ_E, θ_O for the current meta batch:
 Evaluate loss L_{sis} and L_{mim} on D_i^{dt}
 Update θ_O by: $\theta'_O \leftarrow \theta_O - \delta \frac{\sigma L_{total}(f(\theta_E, \theta_O))}{\sigma \theta_O}$
end while

is encoded into one of T possible values, and a classification layer operates within the hidden vector space to indirectly predict the indices of the masks. The loss function is given as:

$$\mathcal{L}_{mim} = - \sum_{k \in M} \log p\left(q_k^\phi(x) | x'\right). \quad (7)$$

Here, $p(q_k^\phi(x) | \tilde{x})$ represents the classification score for classifying the k -th hidden vector belonging to the visual token $q_k^\phi(x)$, where q_ϕ is a categorical distribution.

Due to domain discrepancy, it is difficult to let models learn the intrinsic differences of different domains by themselves. How to mine the key universal information across domains to feedback to the model? To this end, we design the Domain Alignment (DA) loss of each domain and meta-test domain based on the distribution to align the distribution to a specific domain. First, the eigenmeans of the training set are μ_{source} , the covariance matrix is Σ_s , the eigenmeans of the generated samples are μ_s , the covariance matrix is Σ_t .

$$\mathcal{L}_{sis} = \|\mu_s - \mu_t\|^2 + \text{Tr}(\Sigma_s + \Sigma_t - 2(\Sigma_s \Sigma_t)^{1/2}). \quad (8)$$

Based on the feature prior, it is instantiated as they calculate the distance between two distributions with mean and covariance matrices. Smaller distances represent that source domains is closer to the target domain distribution.

3.4 Meta-Domain-Embedding Optimizer

In this subsection, we propose a meta-domain-embedding optimizer based on the MAML [55] paradigm (see Algorithm 1) for pouncing on the generic and personality feature capabilities of learning domain-specific and domain-common features. Here we define each domain as a single task t . In the training process we sample batches of multi-domain data, which consist of meta-train set D_i^{ds} and meta-test set D_i^{dt} , here for simplicity we assume that the full model is described as a function $f(\cdot)$, which receives an image x as input and y as output. The loss function optimized per meta-train domain task during the training is uses cross-entropy loss defined as

$$\begin{aligned} L_{cls}(f(\theta_E, \theta_O)) &= \sum_{(x_j, y_j) \in D_{d_i}} [y_j \log f(x_j) \\ &\quad + (1 - y_j) \log(1 - f(x_j))]. \end{aligned} \quad (9)$$

TABLE 1: Information of the datasets used in our protocols

Source Dataset	Collected from	Synthesis Methods	Identity
FaceForensics++ [14]	YouTube	DeepFake/Face2Face/ FaceSwap/NeuralTextures	-
Celeb-DF (v2) [15]	YouTube	Improved Deepfake	59+
DFDC [17]	Actors	StyleGAN FSGAN Refinement Audioswaps NTH	960
Deepfake in the Wild [16]	Internet	Unknown	100
DeepFakeFace [18]	IMDB/Wikipedia	Stable Diffusion/Inpainting/Insight	-

In this process referred as the inner-loop update, importantly, we just update the learnable token parameter in the meta train and freeze all other feature extraction. θ_E represents the meta-MoE's expert and vpt parameters, while θ_O represents the base model's parameters. After generating the initial domain embeddings θ_i and evaluating the obtained losses on the batch of data, obtains the updated domain embeddings by calculating the gradient of the losses L_{cls} and performing gradient descent updates.

$$\theta'_E \leftarrow \theta_E - \beta \frac{\sigma L_{cls}(f(\theta_E, \theta_O))}{\sigma \theta_E}, \quad (10)$$

where β is the learning rate of gradient descent. In the subsequent step, the model's meta-parameters θ'_E undergo optimization to enhance the performance of meta-test set D_i^{ds} to get the loss L_{cls} and the prediction for domain i .

Similarly, during the meta-test phase, the meta-test sample D_i^{dt} is utilized to update the network. The features are aggregated using the aggregation model after passing through the expert layer. Additionally, the consistency loss L_{sis} of the features is employed to minimize the distance between the source domain and the target domain with reconstructed facial features aid fine-grained forgery feature learning. The overall model loss is stated as follows

$$L_{total} = L_{sis} + L_{cls} + L_{mim}. \quad (11)$$

Then we can optimize the generator $f(\cdot)$ by the gradient:

$$\theta'_O \leftarrow \theta_O - \delta \frac{\sigma L_{total}(f(\theta_E, \theta_O))}{\sigma \theta_O}. \quad (12)$$

In summary, θ_E is updated during the meta-train process to learn the private characteristics of each domain and has higher flexibility due to the dynamic prompt vector. θ_O is updated during the meta-test process to capture generic forged clues, which helps the model acquire complementary information and be used for multi-domain training.

4 MULTI-DOMAIN DEEPPAKE DETECTION PROTOCOLS

Towards the era of large-scale multi-dataset training and cross-dataset testing, we establish a novel benchmark with five mainstream datasets, including FaceForensics++ (FF++) [14], Celeb-DF (v2) (Celeb for short) [15], WildDeepfake (WDF) [16], and DFDC [17]. Information and visualization of these datasets can be found in Table 1 and Figure 5, respectively.

Fig. 5: Visualization of typical samples from five datasets, i.e., FF++ [14], Celeb-DF (v2) [15], DFF [18], WDF [16], and DFDC [17].



4.0.1 FaceForensics++

The FF++ dataset [14] contains video footage of faces that were faked using four common face faking methods: Deepfakes (DF) [1], Face2Face (F2F) [59], FaceSwap (FS) [60], and Nulltextures (NT) [13]. The original video footage was obtained from YouTube, including 1000 real videos and 4000 fake videos. In order to simulate different qualities, the FF++ dataset is available in both high quality (HQ) and low-quality versions (i.e., c23 and c40).

4.0.2 Celeb-DF(V2)

The Celeb-DF(V2) dataset [15] consists of 590 real videos and 5639 fake videos, all of which are 30 seconds long. The original videos come from YouTube public videos and cover a wide distribution of gender, age and ethnicity. Celeb-DF(V2) uses an improved DeepFake algorithm to generate high-resolution faces, which employs a codec with more layers and increased dimensionality. Also, a color conversion algorithm is introduced to address issues such as inconsistent facial colors, and the quality of the generated video is improved by adding training data and post-processing.

4.0.3 DFDC

The DFDC [17] dataset is currently the largest publicly available dataset in the field, containing real videos from 3,426 paid actors. The dataset generates more than 100,000 fake videos through a variety of faking methods, including DeepFakes methods, GAN methods, and non-deep learning methods.

4.0.4 WildDeepfake

This database contains 7,314 facial action sequences extracted from 707 Deepfake videos, all of which are rich and diverse from the web. These facial action sequences are extracted to make the visual effects more realistic and more in line with real-life scenarios.

4.0.5 DFF

A total of 30,000 real images and 90,000 fake images were generated from the original IMDB-WIKI [61] dataset using the Stable Diffusion Inpainting and InsightFace toolbox methods respectively.

Protocols. Although some existing studies have proposed different forgery methods for single-dataset training. No pilot study is available for training on multiple datasets

TABLE 2: Results (AUC (%) and ACC (%)) of joint training on FF++ [14], Celeb-DF [15], and DFF [18] datasets.

Source Domain	Baseline Method	Cross-Domain						In-Domain			
		Tested on DFDC [17]		Tested on WDF [16]		Test On FF++ [14]		Test On Celeb [15]		Test On DFF [18]	
		AUC(%)	ACC(%)	AUC(%)	ACC(%)	AUC(%)	ACC(%)	AUC(%)	ACC(%)	AUC(%)	ACC(%)
FF++ [14]	Xception [13] (<i>ICCV</i> 2019)	75.49	53.82	62.74	57.36	100	97.87	89.10	90.24	89.73	93.12
	REECE [12] (<i>CVPR</i> 2022)	75.19	73.42	77.90	62.18	99.98	98.17	92.31	94.16	93.11	94.38
	UCF [46] (<i>ICCV</i> 2023)	80.50	73.01	73.40	67.52	99.60	99.60	82.40	86.14	93.11	94.38
	Implicit [12] (<i>CVPR</i> 2023)	74.90	72.13	75.12	69.40	99.98	96.23	82.80	83.19	90.11	92.80
	CLIP [29] (<i>ICML</i> 2021)	76.01	72.51	74.33	64.52	93.21	96.10	81.43	83.71	92.21	93.19
Celeb [15]	Xception [13] (<i>ICCV</i> 2019)	53.12	52.16	66.67	43.75	51.32	54.62	96.32	98.61	73.10	75.46
	REECE [12] (<i>CVPR</i> 2022)	57.26	54.71	69.32	67.15	53.17	55.71	99.20	99.33	76.32	74.53
	UCF [46] (<i>ICCV</i> 2023)	65.04	62.90	61.24	63.90	61.46	63.09	97.12	97.06	73.71	72.33
	Implicit [56] (<i>CVPR</i> 2023)	64.60	62.12	65.11	64.54	61.08	65.71	99.20	93.33	76.32	71.02
	CLIP [29] (<i>ICML</i> 2021)	54.32	51.67	65.17	61.34	51.03	52.00	96.98	93.12	72.33	76.45
FF++ [14] & Celeb [15]	Xception [13] (<i>ICCV</i> 2019)	67.19	56.62	59.31	56.22	82.99	68.75	100	96.87	93.89	91.23
	REECE [12] (<i>CVPR</i> 2022)	70.32	63.18	64.61	62.83	85.24	73.75	100	98.25	94.21	93.54
	UCF [46] (<i>ICCV</i> 2023)	65.43	63.21	67.50	63.96	85.72	82.10	97.98	97.04	89.32	87.49
	Implicit [56] (<i>CVPR</i> 2023)	57.26	67.81	63.50	81.70	78.50	75.41	98.41	98.10	89.63	86.10
	CLIP [29] (<i>ICML</i> 2021)	67.41	62.78	65.78	63.08	83.51	76.10	100	97.71	93.53	92.15
FF++ [14] & Celeb [15] & DFF [18]	Xception [13] (<i>ICCV</i> 2019)	58.82	53.12	58.82	53.12	96.13	89.26	99.79	96.55	93.93	95.41
	REECE [12] (<i>CVPR</i> 2022)	73.16	67.09	66.04	63.00	97.12	91.19	99.38	98.19	95.88	96.17
	UCF [46] (<i>ICCV</i> 2023)	67.40	59.70	72.10	69.54	85.29	83.50	98.74	97.41	92.32	92.08
	Implicit [56] (<i>CVPR</i> 2023)	68.91	61.43	69.32	69.54	89.32	89.40	99.20	99.33	90.04	90.56
	CLIP [29] (<i>ICML</i> 2021)	55.43	51.02	73.45	71.07	96.26	85.95	99.24	97.23	94.71	93.23
	GM-DF (Ours)	77.54	75.23	79.70	75.08	98.23	97.23	99.99	98.45	97.72	98.78

TABLE 3: The results on the Multi-Domain Deepfake detection benchmarks based on $M_{EER}(\%)$ and AUC (%).

Method	FF++ & Celeb & DFF		FF++ & Celeb & WDF		FF++ & DFF & WDF		Celeb & DFF & WDF	
	$M_{EER}(\%)$	AUC(%)	$M_{EER}(\%)$	AUC(%)	$M_{EER}(\%)$	AUC(%)	$M_{EER}(\%)$	AUC(%)
MesoNet [57] (<i>WIFS</i> 2018)	45.31	53.34	44.70	69.25	46.42	57.16	40.54	54.92
Multi-task [58] (<i>BTAS</i> 2019)	36.33	57.33	36.98	66.03	39.14	74.72	34.91	62.17
F ³ -Net [10] (<i>ECCV</i> 2020)	36.12	57.76	35.82	68.35	37.23	67.53	31.05	66.89
Xception [13] (<i>ICCV</i> 2021)	33.45	71.09	37.68	66.64	35.56	76.88	33.30	65.40
REECE [12] (<i>CVPR</i> 2022)	32.57	70.85	35.07	69.86	36.64	78.72	30.14	71.92
UCF [46] (<i>ICCV</i> 2023)	35.90	69.72	34.78	65.41	35.02	74.13	34.02	69.11
Implicit [56] (<i>CVPR</i> 2023)	33.09	69.54	37.66	72.12	38.91	74.10	33.18	68.31
GM-DF (Ours)	30.13	74.33	32.81	72.37	32.90	80.19	28.36	73.73

with real-world diverse forgery patterns and large-scale characteristics. Besides the perspective of training, only a single forgery test domain is usually used in the evaluation of algorithm performance, which leads to biased comparisons of state-of-the-art methods. To tackle the above-mentioned issues, we provide a novel data arrangement and training/testing strategy to benchmark the fair evaluations. Specifically, 5 datasets (each dataset is regarded as an individual domain), i.e., FF++ [14], WDF [16], Celeb [15], DFDC [17], and DFF [18] are merged into a large set D , which can be further divided into training sets $\{D_{FF++}, D_{WDF}, D_{Celeb}, D_{DFF}\}$ and test set $\{D_{DFDC}, D_i\}$. $i \in \{FF++, WDF, Celeb, DFF\}$ which denotes the subsets removed from the training set for the testing set. Specifically, in consideration of costly training time, the large-scale DFDC [17] is only used for testing. We randomly select $n \leq 3$ subsets of the data for training. $n = 1$ for the traditional single-domain training protocol while $n = 3$ denotes the newly established multi-domain protocols: $\{D_{FF++} \cup D_{WDF} \cup D_{Celeb}\}$, $\{D_{FF++} \cup D_{WDF} \cup D_{DFF}\}$, $\{D_{FF++} \cup D_{DFF} \cup D_{DFF}\}$ and $\{D_{WDF} \cup D_{Celeb} \cup D_{DFF}\}$ for training, respectively. More details can be found in the supplementary material.

Evaluation metrics. Three common metrics, i.e., Accuracy (ACC (%)), Area Under ROC Curve (AUC (%)), and Equal Error Rate (EER (%)) are adopted. EER is defined as the error rate when false acceptance rate (FAR) is equal to false rejection rate (FRR), and can be expressed as $EER = \frac{FRR + FAR}{2}$. However, in the realm of evaluating face forgery detection, we are more concerned with keeping the FAR at a relatively low level to ensure that it will not be easy to authenticate a

forged face. For this purpose, we introduce the a priori probability of positive examples when calculating the EER. Since the impact on the system of a positive sample misclassified as a negative example is much greater than the impact of a negative sample as a positive example. To counteract this effect, we introduced the P_{real} parameter. In addition, we found that the original EER did not take into account the effect of testing on multiple domains. Calculating the EER directly on each dataset and then averaging the values may be affected by extreme values, and we judged performance against multiple domains by taking the maximum of instead of simply averaging them. This ensures that our evaluation is more accurate and robust.

$$M^i = P_{real}^i * FRR^i + (1 - P_{real}^i) * FAR^i$$

$$M_{ERR} = \text{Max} \{M^1, M^2, \dots, M^N\} \quad (13)$$

Here P_{real} is the prior probability of the real samples. M^i is the prior probability EER of the i -th target domain. N is the number of target test domains.

5 EXPERIMENTS

In this section, we evaluate the performance of our proposed method on FaceForensics++ (FF++) [14], Celeb-DF(V2) [15], DFDC [17] and WildDeepfake [16] datasets on both traditional protocols as well as the proposed benchmark.

5.1 Implementation Details

We use ViT-B/16 [29] as the backbone model. We use RetinaFace [62] to detect facial areas and scaled the face

image to 224×224 with a patch size of 16. We trained the model using the Adam optimizer with the learning rate set to $3e-6$. The batch size during training was 32, and 40 training epochs were performed. Following official setting [14], we extracted 100 frames from each video as validation set and test set. During the training process, only random flipping was used for data augmentation. During the dataset merging phase, multiple datasets are randomly shuffled and then consolidated into a new dataset. All code is implemented using the PyTorch framework.

5.2 Preliminary Multi-datasets Investigation

Single-dataset: Each dataset is trained separately to be evaluated on different test sets, and the cross-dataset performance of the model is tested on different dataset datasets using the model trained from scratch.

Direct-Merged datasets: After directly merging multiple face forgery datasets, we employ a straightforward strategy of training a baseline fake detector using a generalized classification loss. This aims to verify whether directly merging the datasets is expected to improve the existing forgery detection models across datasets.

To investigate the feasibility of co-training from multiple datasets to improve the cross-dataset performance, we use multiple datasets from different sources for training. We also add the recently proposed which use stable diffusion generated face data for exploration, we believe is more in line with the real-world nature of forgery methods. Table 2 shows the following essential findings:

1) *Directly combining datasets for training does not increase accuracy.* We first train the baseline detector on a single dataset (e.g., FF++ [14] or Celeb [15] as well as DFF [18]) and evaluate this trained baseline on two different datasets (eg. WDF [16] and DFDC [17]). As shown in Table 2 the baseline performs well only on its original training dataset (e.g., 75.49% AUC to 67.19%). Its detection accuracy drops severely when the baseline detector combines the dataset Celeb-DF(V2) (with an accuracy of only 53.12%). This is mainly due to the fact that the single dataset detection model over-fits the common features of its training dataset, but does not take into account the variations in the characteristics of the source-to-target dataset. The same accuracy degradation problem can be observed on other baseline detectors such as REECE [12] and recent sota model [46], [56].

2) *Pre-trained model performs well within domain and modelling knowledge fading in unseen samples.* We fine-tuned the baseline model on both individual and multiple datasets, and subsequently compared the outcomes. We observed that the model trained on either the FF++ or Celeb dataset performed better when trained individually, whereas its performance deteriorated when trained jointly on both datasets. Table 2 shows that the model is pretrained on FF++ and Celeb-DF(V2), but the detection accuracy of DFDC [17] is still poor because the model has been fine-tuned to Celeb [15], forgetting what was learned from the previous pretraining dataset and the differences in joint training at different resolutions.

3) *Training stable diffusion and GAN face forgery data together may increase these differences and learning difficulties.*

For example, after fusing DFF [18] (the dataset generated by the diffusion model) in the WDF [16] test the results of the Xception [13] model dropped the most from 67.19% to 58.82%, which has a higher demand on the detector. Both recent sota model UCF [46] and Implicit [56] exhibit similar characteristics in the experimental results, leading to an unavoidable performance decline. Specifically, UCF [46] achieves 97.10% on Celeb training and testing, but rapidly drops to 73.71% when tested on the DFF dataset.

4) *The effectiveness of the proposed method.* Table 2 shows that the average cross- and intra-domain detection performance of the proposed method exceeds the baseline models, proving its flexibility and validity. Our model demonstrates an average AUC improvement of 8.87% over UCF [46] in cross-domain scenarios and 6.53% in within-domain scenarios, underscoring the superiority of the two-stage learning approach.

5.3 Results on the Proposed Protocols

To further assess the real-world performance of our method, we conducted experiments on the proposed multi-datasets deepfake detection benchmark (in Sec. 4). We compared our method with commonly used forgery detection networks such as MesoNet [57] and Xception [13], as well as some recent state-of-the-art (SOTA) methods. RECCE [12], UCF [46] and Implicit [56] was tested under default settings. As shown in Table 3, according to the results under evaluation metrics M_{ERR} & AUC, the proposed method outperforms other methods, showcasing its effectiveness. An intriguing finding emerged: Despite achieving excellent performance on individual datasets, some existing SOTA methods experience drastic performance drop under multi-datasets protocols.

5.4 Results on Traditional Protocols

We also conducted experiments on the commonly used benchmarks [69] to prove the effectiveness of our methodology in classic single dataset multi-source operation settings. For this experimental setup, we selected one class of manipulated forged videos from FF++ [14] as the unseen manipulation sample, while utilizing the remaining three classes as the training set. The experiment results are in Table 4. All four assessment settings show our model achieved better detection results. Our Model achieves 3.13% improvement in ACC on GID-DF (C23) compared to recent sota model Implicit [56], proving that our model can better adapt to forgery methods in multiple source domains and learn the common and characteristic features of each forgery method.

To gain a more detailed understanding of cross-domain performance, we employed a single data training approach. Specifically, our model was trained on FF++ (C23) [14] and tested on DFDC [17] and Celeb-DF(V2) [15]. Comparative results with state-of-the-art methods are presented in Table 5. In cross-dataset comparisons, our model demonstrates excellent performance. In internal tests, when compared to the recent REECE [12], our model exhibits a notable improvement of 4.17% and 3.59% compare to SFGD [76] in AUC. This indicates that our model outperforms traditional cross-forgery patterns and other state-of-the-art models in terms of generalization, showcasing the transfer capabilities

TABLE 4: Cross-Manipulation Evaluation: ACC (%) and AUC (%) for Multi-Source Training and Testing.

Method	GID-DF (C23)		GID-DF (C40)		GID-F2F (C23)		GID-F2F (C40)	
	ACC(%)	AUC(%)	ACC(%)	AUC(%)	ACC(%)	AUC(%)	ACC(%)	AUC(%)
EfficientNet [63] (<i>PMLR</i> 2019)	82.40	91.11	67.60	75.30	63.32	80.10	61.41	67.40
ForensicTransfer [64] (<i>Arxiv</i> 2018)	72.01	-	68.20	-	64.50	-	55.00	-
Multi-task [65] (<i>BTAS</i> 2019)	70.30	-	66.76	-	58.74	-	56.50	-
F ³ -Net [66] (<i>ECCV</i> 2020)	83.57	94.95	77.50	85.77	61.07	81.20	64.64	73.70
MLGD [67] (<i>AAAI</i> 2018)	84.21	91.82	67.15	73.12	63.46	77.10	58.12	61.70
LTW [68] (<i>AAAI</i> 2021)	85.60	92.70	69.15	75.60	65.60	80.20	65.70	72.40
DCL [69] (<i>AAAI</i> 2022)	87.70	94.90	75.90	83.82	68.40	82.93	67.85	75.07
M2TR [70] (<i>ICML</i> 2022)	81.07	94.91	74.29	84.85	55.71	76.99	66.43	71.70
Implicit [56] (<i>CVPR</i> 2023)	88.21	95.03	76.90	84.55	69.36	84.37	67.99	74.80
GM-DF (Ours)	91.34	96.62	78.13	85.19	72.32	86.30	69.02	75.51

TABLE 5: Cross-domain comparisons of generalization based on AUC (%). We train the model on the HQ dataset of FF++ [14] and evaluate it on Celeb-DF(V2) [15] and DFDC [17].

Method	Celeb	DFDC
EN-B4 [71] (<i>PMLR</i> 2019)	66.24	66.81
F ³ -Net [10] (<i>ECCV</i> 2020)	71.21	72.88
Xception [13] (<i>ICCV</i> 2021)	66.91	69.93
MAT(EN-B4) [72] (<i>CVPR</i> 2021)	76.65	67.34
Face X-ray [22] (<i>CVPR</i> 2021)	74.20	70.00
RFM [73] (<i>CVPR</i> 2021)	67.64	68.01
SRM [8] (<i>CVPR</i> 2021)	79.40	79.70
Local-relation [74] (<i>AAAI</i> 2021)	78.26	76.53
LTW [75] (<i>AAAI</i> 2021)	77.14	74.58
RECCE [12] (<i>CVPR</i> 2022)	77.39	76.75
Impliciry [56] (<i>ICCV</i> 2023)	82.04	-
SFGD [76] (<i>CVPR</i> 2023)	75.83	73.64
GM-DF (Ours)	83.16	77.23

TABLE 6: Natural language descriptions of the real and fake face used to train the model. BLIP Generate indicates that the BLIP [77] model generates descriptive information.

Prompt	Real Prompt	Fake Prompt
P1	A photo of real face	A photo of fake face
P2	This is a photo of real	This is a photo of fake
P3	{BLIP Generate} A photo of real face	{BLIP Generate} A photo of real face
P4	Real	Fake
P5	This is how a real face looks like	This is how a fake face looks like
P6	This photo contains real face	This photo contains fake face
P7	Real face is in this photo	Fake face is in this photo

of GM-DF models. This underscores the effectiveness of natural language supervision in generating more generalizable representations, particularly in the context of cross-dataset training data.

5.5 Ablation Study

In this subsection, we train on Celebdf [15] and FF++ [14] datasets and test on WDF [16] to test our proposed module and essential parameter settings.

Effectiveness of different text prompts. To validate the effect of different prompts on experimental performance, we introduced new templates into the prompts group. In Table 6, shows the specific language descriptions of the real and fake face categories. we scrutinize the impact of distinct

TABLE 7: Impact of different text prompts (described in Table 6).

Method	FF + Celeb->WDF		FF +DFE->WDF		Celeb +DFE->WDF	
	AUC(%)	EER(%)	AUC(%)	EER(%)	AUC(%)	EER(%)
P1	63.08	29.23	76.09	30.95	78.09	31.27
P2	61.30	29.85	75.88	34.56	77.92	31.92
P3	61.37	31.81	69.93	32.03	74.25	34.87
P4	62.11	30.67	70.25	31.87	72.22	33.89
P5	63.33	30.35	74.40	33.30	77.18	31.46
P6	60.43	35.17	72.20	34.51	76.54	34.46
P7	61.10	33.29	72.43	34.69	78.12	33.15

TABLE 8: Ablation of each component on the protocol of FF++& Celeb& DFF to WDF.

ID	Baseline	DA	MIM	Meta-MoE	AUC	ACC
num1	✓				73.45	71.07
num2	✓	✓			74.36	71.65
num3	✓		✓		75.11	73.33
num4	✓			✓	77.21	74.18
num5	✓	✓	✓		75.71	73.42
num6	✓	✓		✓	75.12	74.39
GM-DF (Ours)	✓	✓	✓	✓	79.70	75.13

text prompts on the model. Notably, varied texts exhibit commendable performance across diverse datasets, with marginal differences. This substantiates the notion that text can effectively manifest dynamicized parameters in real-world contexts, thereby affirming our concept of instating dynamic affine transformations tailored to each dataset. An intriguing discovery emerges when utilizing BLIP [77] to generate images with detailed descriptive information alongside the original combination of category images. Surprisingly, performance experiences relative degradation, potentially attributed to interference induced by category-independent prompts.

Impacts of various ViT backbone initialization. To extend our observations on the impact of initialization on the multi-datasets training, we tuned the model using different CLIP pre-training weights and showed a comparison of their performance in Table 7. Specifically, we fine-tuned the weighting using two architectures, Resnet and Vit, a) Resnet50 backbone; b) Resnet101 backbone; c) ViT backbone with a patch size of 16; d) ViT backbone with a patch size of 32; and e) ViT backbone with a patch size of 14. It can be seen that ViT pre-training initialization yields better multi-dataset training generalization compared to other initialization methods. Compared to other initialization methods, the Transformer initialization achieves better multi-datasets training generality due to its powerful representation extraction capability, which provides a better image-text alignment basis and detailed feature extraction capability for all image alignment experiments.

Effectiveness of DA loss. In the 8 first and second rows compared to the baseline, the DA loss achieved an improvement of approximately 0.91%, demonstrating the need for the alignment of the distributions of the two feature datasets through higher-order statistical features. We can observe a consistent improvement in performance when using the DA loss function, which demonstrates the advantage of dataset alignment with the visual-linguistic pre-training model.

Effectiveness of MIM loss. Comparing the first and third

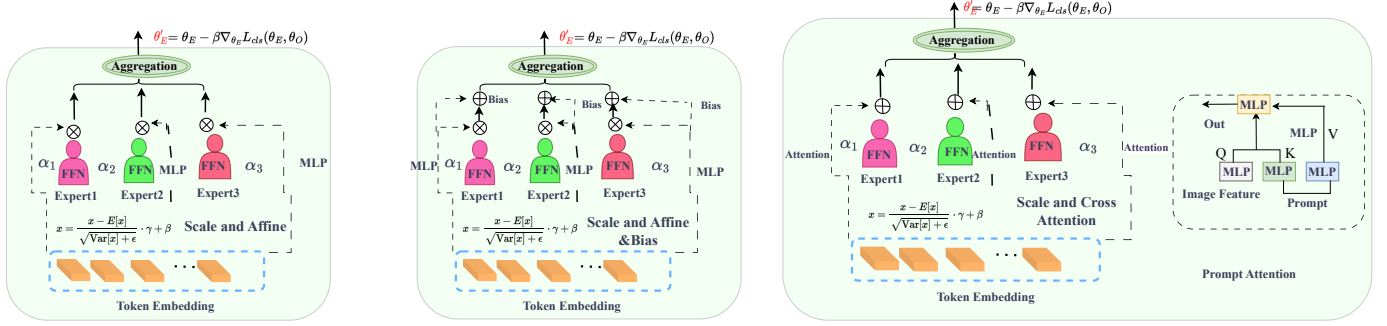


Fig. 6: Architectures of three adaptation strategies for the Dataset Information Layer, including Affine (left), Affine&Bias (middle), and Cross Attention (right).

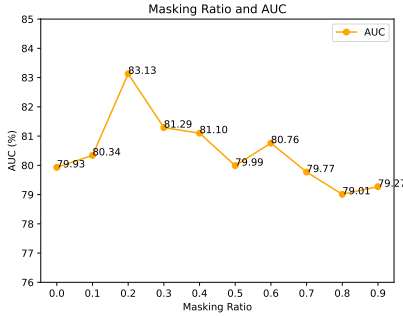


Fig. 7: Quantitative analyses of masking strategy. The AUC (%) scores of cross-dataset evaluation on Celeb-DF are reported.

rows, it can be seen that the addition of the reconstruction module improves the AUC by 1.66% over the original model, which indicates that the reconstructed features can effectively enhance the ability of fine-grained information extraction on the forged face.

TABLE 9: Results of different domain adaptive strategies when trained on Celeb-DF (v2) [15] & DFF [18] and tested on WDF [16] dataset.

Method	Celeb-DF (v2) & DFF	
	ACC(%)	AUC(%)
Affine	65.21	67.32
Affine & Bias	63.24	66.13
Cross Attention	64.18	66.87

Effectiveness of Meta-MoE. To quantify the importance of Meta-Moe module, we compare our text-based supervisory signals with meta learning and without two stage learning. It can be seen from the fourth line that meta-MoE plays an important role in performance improvement (from 73.45% to 77.21%), which is mainly caused by learning the characteristic features of the domain. The mask-supervision method exhibits better generalization, suggesting that mask supervision alone can restrain overfitting to the training data. Moreover, unlike the only text backbone we improves steadily with more fine-grained supervision,, which further confirms the scalability and versatility of multidataset learning .

Analysis of masking ratio. The quantitative results of the

cross-dataset evaluation are shown in Figure 7. We observe that the minimum and randomized masking strategy achieves optimality under medium masking rates. Their performance is severely degraded as the masking rate is greater than 80%. The random masking strategies work best at 20% maskingrate. This indicates that some important face edges may be corrupted using the random masking strategy.

5.6 Visualization and Discussion

Discussion about the Dataset Information Layer. To address the challenge of feature adaptation to different dataset domains, as illustrated in Figure 5, we also investigate three different domain adaptation strategies (i.e., Affine, Affine&Bias, and Cross Attention) for the Dataset Information Layer.

1) Affine. The domain-specific knowledge of each domain can be realized by linearly mapping the respective domain prompt feature multiplication to the intermediate feature layer. This part of the linear mapping is realized through a single MLP. We also visualized the Adaptive layer features, which suggests that the differences of different data sets are effectively learned.

2) Affine&Bias. Here, we adjust the LayerNorm’s parameters via learning by both affines and offsets. The vanilla LayerNorm assumes that the samples are all from the same distribution, which the data might come from different domains. Therefore, the parameters in LayerNorm should be not the same in different domains. The LayerNorm based on Affine&Bias learning can be formulated as follows:

$$\text{LayerNorm}(x) = \frac{x - E[x]}{\sqrt{\text{Var}[x] + \epsilon}} \cdot (\gamma * \gamma_d) + (\beta_d). \quad (14)$$

Domain-specific parameters $\gamma * \gamma_d$ and $\beta + \beta_d$ can adaptively change the intermediate representation conditions and domain indicators capture distinctive characteristics.

3) Cross Attention. Based on [42], [78], we use Cross Attention to aggregate text features and raw features with a cross attention block with jump connections at the beginning of each encoder-decoder stage. First, we partitioned the domain cues into n independent in-domain cue embeddings that have the same shape, which partially acts as a reference set for cross-attention, with the images providing the associated information. Next, a series of attention operations are performed between the query vectors generated for each image and the key-value vectors generated for the domain

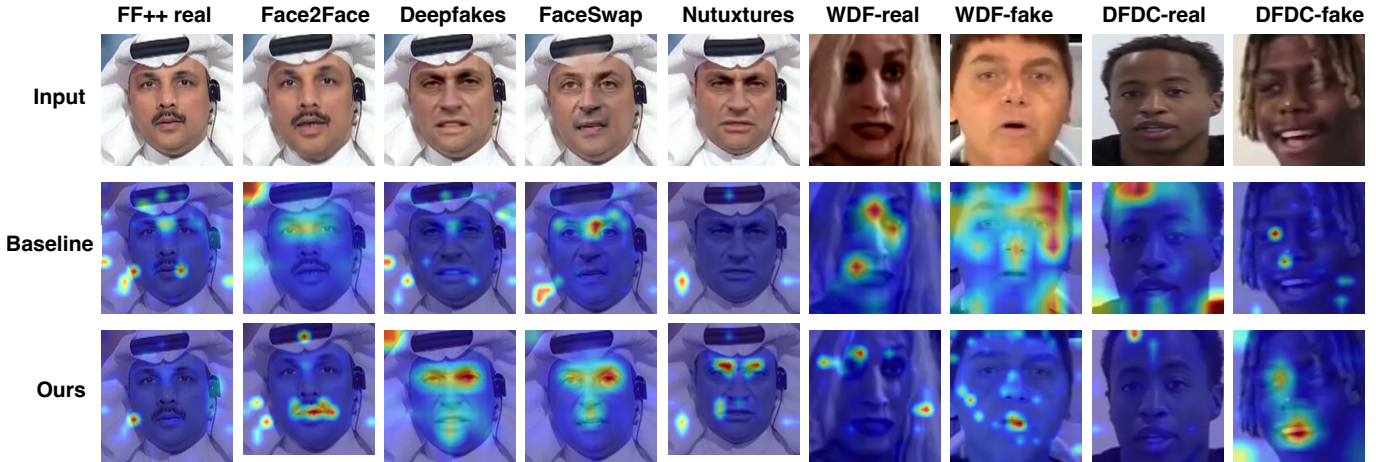


Fig. 8: The model’s attention is illustrated through a heatmap, where darker colors signify increased focus in that specific region. The first column represents the input image, the second column depicts the outcome obtained by directly fusing the data using fine-tuning with the CLIP [29], and the third column showcases the results from our model.

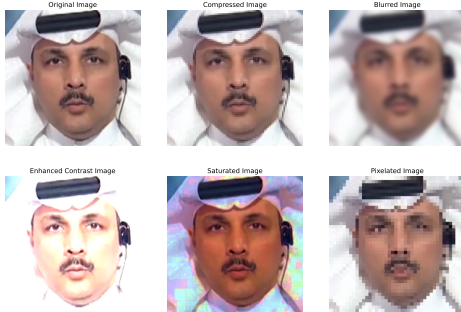


Fig. 9: Examples of images with different quality-degradation methods. Image Compression, Gaussian Blur, Enhanced Contrast dithering, Saturated dithering, and pixelization, respectively.

cues. Finally, the results of the attention operations are added to the data point embeddings after projection by a zero-initialized linear layer. To validate the effectiveness of our model in Figure 4 we use to visualize the ROC curves, the data were trained in FF++c23 and Celeb and tested on various datasets .

The results of these three domain adaptive strategies when trained on Celeb-DF (v2) [15] & DFF [18] and tested on WDF [16] are shown in Table 9. We can find that the Affine strategy is simple yet effective, and achieves better cross-domain performance than other two alternatives. Besides, we also find that the performance of Cross Attention strategy seems satisfactory, and one possible future direction is how to efficiently combine Affine with Cross Attention to boost generalization capacity.

Analysis of robustness against distortions. Considering the prevalence of image processing on the web, we investigate the performance under several distortions proposed by [12], [58], namely image compression, Gaussian blurring, contrast dithering, saturation dithering and pixelization. The quality-degraded images using different degradation methods are shown in Fig 8. The results are shown in Table 10. We can see that our model is more robust to the listed

TABLE 10: Robustness evaluation in terms of AUC (%) on WildDeepfake (WDF) dataset.

Method	Compress	Blur	Contrast	Saturate	Pixelate	Avg
Multi-task [58] (<i>BTAS</i> 2019)	89.64	80.98	89.30	90.37	79.44	85.95
F ³ -Net [10] (<i>ECCV</i> 2020)	86.71	78.99	86.53	87.67	73.23	82.63
Xception [13] (<i>ICCV</i> 2021)	86.01	78.29	81.90	84.96	66.24	79.48
RFM [73] (<i>ECCV</i> 2021)	83.74	75.34	79.77	82.59	71.25	78.54
Add-Net [79] (<i>AAAI</i> 2021)	83.34	79.66	84.46	85.13	64.33	79.38
REECE [12] (<i>CVPR</i> 2022)	89.65	87.29	91.19	91.74	83.88	88.75
GM-DF (Ours)	90.32	89.43	92.56	92.31	84.95	89.91

ingressions than the existing methods. Both our method and previous methods are generally robust to compression, contrast and saturation. However, in scenarios blur and pixelate, the performance [10], [12], [13], [58], [73] are still much lower than the proposed method, indicating the robustness of the proposed method.

Visualization. We employed a joint training approach using three datasets FF++ [14], Celeb-DF (V2) [15], and DFF [18]. Subsequently, we conducted visual analyses on individual in-domain datasets as well as various cross-domain datasets. From Figure 8, it can be observed that directly merging datasets often leads the model to lose effective focus in challenging scenarios, such as WDF [16], where attention shifts to background regions. In contrast, our proposed multi-domain fusion model consistently concentrates on facial regions and successfully detects manipulated faces.

6 CONCLUSION

In this paper, we investigate the generalization capacity of deepfake detectors when trained on multi-dataset scenarios and propose a novel benchmark for multi-scenario training. We design a Generalized Multi-Scenario Deepfake Detection (GM-DF) framework to learn of both specific and common features across datasets. By utilizing generic text representations to learn the relationships across different datasets, we propose a novel meta-learning strategy to capture the relational information among datasets. Besides, GM-DF employs contrastive learning on image-text pairs to capture common dataset characteristics and utilizes self-supervised mask relation learning to mask out partial cor-

relations between regions during training. Extensive experiments demonstrate the superior generalization of our method. In the future, we plan to explore techniques for localizing counterfeit regions and enhancing generalization by leveraging multimodal large language models.

REFERENCES

- [1] "Deepfakes," <https://github.com/deepfakes/faceswap>, accessed 2022-10-29.
- [2] "Faceswap," <https://github.com/MarekKowalski/FaceSwap>, accessed 2022-10-29.
- [3] J. Thies, M. Zollhofer, M. Stamminger, C. Theobalt, and M. Nießner, "Face2face: Real-time face capture and reenactment of rgb videos," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2387–2395.
- [4] J. Thies, M. Zollhöfer, and M. Nießner, "Deferred neural rendering: Image synthesis using neural textures," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 4, pp. 1–12, 2019.
- [5] Z. Yu, R. Cai, Z. Li, W. Yang, J. Shi, and A. C. Kot, "Benchmarking joint face spoofing and forgery detection with visual and physiological cues," *IEEE Transactions on Dependable and Secure Computing*, 2024.
- [6] Y. Shi, Y. Gao, Y. Lai, H. Wang, J. Feng, L. He, J. Wan, C. Chen, Z. Yu, and X. Cao, "Shield: An evaluation benchmark for face spoofing and forgery detection with multimodal large language models," *arXiv preprint arXiv:2402.04178*, 2024.
- [7] H. H. Nguyen, J. Yamagishi, and I. Echizen, "Capsule-forensics: Using capsule networks to detect forged images and videos," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 2307–2311.
- [8] Y. Luo, Y. Zhang, J. Yan, and W. Liu, "Generalizing face forgery detection with high-frequency features," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 16 317–16 326.
- [9] J. Frank, T. Eisenhofer, L. Schönherr, A. Fischer, D. Kolossa, and T. Holz, "Leveraging frequency analysis for deep fake image recognition," in *International conference on machine learning*. PMLR, 2020, pp. 3247–3258.
- [10] Y. Qian, G. Yin, L. Sheng, Z. Chen, and J. Shao, "Thinking in frequency: Face forgery detection by mining frequency-aware clues," in *European conference on computer vision*. Springer, 2020, pp. 86–103.
- [11] Z. Chen, L. Xie, S. Pang, Y. He, and B. Zhang, "Magdr: Mask-guided detection and reconstruction for defending deepfakes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9014–9023.
- [12] J. Cao, C. Ma, T. Yao, S. Chen, S. Ding, and X. Yang, "End-to-end reconstruction-classification learning for face forgery detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 4113–4122.
- [13] A. Rossler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner, "Faceforensics++: Learning to detect manipulated facial images," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 1–11.
- [14] —, "Faceforensics++: Learning to detect manipulated facial images," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 1–11.
- [15] Y. Li, X. Yang, P. Sun, H. Qi, and S. Lyu, "Celeb-df: A large-scale challenging dataset for deepfake forensics," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 3207–3216.
- [16] B. Zi, M. Chang, J. Chen, X. Ma, and Y.-G. Jiang, "Wilddeepfake: A challenging real-world dataset for deepfake detection," in *Proceedings of the 28th ACM international conference on multimedia*, 2020, pp. 2382–2390.
- [17] B. Dolhansky, J. Bitton, B. Pflaum, J. Lu, R. Howes, M. Wang, and C. C. Ferrer, "The deepfake detection challenge (dfdc) dataset," *arXiv preprint arXiv:2006.07397*, 2020.
- [18] H. Song, S. Huang, Y. Dong, and W.-W. Tu, "Robustness and generalizability of deepfake detection: A study with diffusion models," 2023.
- [19] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [20] A. A. Pokroy and A. D. Egorov, "Efficientnets for deepfake detection: Comparison of pretrained models," in *2021 IEEE conference of russian young researchers in electrical and electronic engineering (ElConRus)*. IEEE, 2021, pp. 598–600.
- [21] D. Güera and E. J. Delp, "Deepfake video detection using recurrent neural networks," in *2018 15th IEEE international conference on advanced video and signal based surveillance (AVSS)*. IEEE, 2018, pp. 1–6.
- [22] H. Liu, X. Li, W. Zhou, Y. Chen, Y. He, H. Xue, W. Zhang, and N. Yu, "Spatial-phase shallow learning: rethinking face forgery detection in frequency domain," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 772–781.
- [23] L. Li, J. Bao, T. Zhang, H. Yang, D. Chen, F. Wen, and B. Guo, "Face x-ray for more general face forgery detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 5001–5010.
- [24] Y. Lai, Z. Luo, and Z. Yu, "Detect any deepfakes: Segment anything meets face forgery detection and localization," in *Chinese Conference on Biometric Recognition*, 2023, pp. 180–190.
- [25] T. Zhao, X. Xu, M. Xu, H. Ding, Y. Xiong, and W. Xia, "Learning self-consistency for deepfake detection," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 15 023–15 033.
- [26] L. Chen, Y. Zhang, Y. Song, L. Liu, and J. Wang, "Self-supervised learning of adversarial example: Towards good generalizations for deepfake detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 18 710–18 719.
- [27] C. Kong, K. Zheng, Y. Liu, S. Wang, A. Rocha, and H. Li, " m^3 fas: An accurate and robust multimodal mobile face anti-spoofing system," *IEEE Transactions on Dependable and Secure Computing*, 2024.
- [28] R. Zhang, W. Zhang, R. Fang, P. Gao, K. Li, J. Dai, Y. Qiao, and H. Li, "Tip-adapter: Training-free adaption of clip for few-shot classification," in *European Conference on Computer Vision*. Springer, 2022, pp. 493–510.
- [29] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [30] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Conditional prompt learning for vision-language models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 816–16 825.
- [31] X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui, "Open-vocabulary object detection via vision and language knowledge distillation," *arXiv preprint arXiv:2104.13921*, 2021.
- [32] X. Zhai, X. Wang, B. Mustafa, A. Steiner, D. Keysers, A. Kolesnikov, and L. Beyler, "Lit: Zero-shot transfer with locked-image text tuning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 123–18 133.
- [33] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds *et al.*, "Flamingo: a visual language model for few-shot learning," *Advances in Neural Information Processing Systems*, vol. 35, pp. 23 716–23 736, 2022.
- [34] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," *Advances in neural information processing systems*, vol. 36, 2024.
- [35] J. Li, D. Li, S. Savarese, and S. Hoi, "Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models," *arXiv preprint arXiv:2301.12597*, 2023.
- [36] R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.
- [37] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.
- [38] H. Huang, L. Lin, R. Tong, H. Hu, Q. Zhang, Y. Iwamoto, X. Han, Y.-W. Chen, and J. Wu, "Unet 3+: A full-scale connected unet for medical image segmentation," in *ICASSP 2020-2020 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2020, pp. 1055–1059.
- [39] H. Cao, Y. Wang, J. Chen, D. Jiang, X. Zhang, Q. Tian, and M. Wang, "Swin-unet: Unet-like pure transformer for medical image segmentation," in *European conference on computer vision*. Springer, 2022, pp. 205–218.
- [40] X. Dai, Y. Chen, B. Xiao, D. Chen, M. Liu, L. Yuan, and L. Zhang, "Dynamic head: Unifying object detection heads with attentions,"

- in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 7373–7382.
- [41] R. Gong, D. Dai, Y. Chen, W. Li, and L. Van Gool, “mdalu: Multi-source domain adaptation and label unification with partial datasets,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 8876–8885.
- [42] X. Wang, Z. Cai, D. Gao, and N. Vasconcelos, “Towards universal object detection by domain attention,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 7289–7298.
- [43] X. Zhou, V. Koltun, and P. Krähenbühl, “Simple multi-dataset detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 7571–7580.
- [44] J. Lambert, Z. Liu, O. Sener, J. Hays, and V. Koltun, “Mseg: A composite dataset for multi-domain semantic segmentation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2879–2888.
- [45] X. Zhao, S. Schuster, G. Sharma, Y.-H. Tsai, M. Chandraker, and Y. Wu, “Object detection with a unified label space from multiple datasets,” in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*. Springer, 2020, pp. 178–193.
- [46] Z. Yan, Y. Zhang, Y. Fan, and B. Wu, “Ucf: Uncovering common features for generalizable deepfake detection,” *arXiv preprint arXiv:2304.13949*, 2023.
- [47] M. Zhang, H. Marklund, N. Dhawan, A. Gupta, S. Levine, and C. Finn, “Adaptive risk minimization: Learning to adapt to domain shift,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 23 664–23 678, 2021.
- [48] J. Wang, Z. Wu, W. Ouyang, X. Han, J. Chen, Y.-G. Jiang, and S.-N. Li, “M2tr: Multi-modal multi-scale transformers for deepfake detection,” in *Proceedings of the 2022 international conference on multimedia retrieval*, 2022, pp. 615–623.
- [49] S. Masoudnia and R. Ebrahimpour, “Mixture of experts: a literature survey,” *Artificial Intelligence Review*, vol. 42, pp. 275–293, 2014.
- [50] M. Jia, L. Tang, B.-C. Chen, C. Cardie, S. Belongie, B. Hariharan, and S.-N. Lim, “Visual prompt tuning,” in *European Conference on Computer Vision*. Springer, 2022, pp. 709–727.
- [51] Z. Yang, J. Liang, Y. Xu, X.-Y. Zhang, and R. He, “Masked relation learning for deepfake detection,” *IEEE Transactions on Information Forensics and Security*, 2023.
- [52] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, “Masked autoencoders are scalable vision learners,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 000–16 009.
- [53] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola, “What makes for good views for contrastive learning?” *Advances in Neural Information Processing Systems*, vol. 33, pp. 6827–6839, 2020.
- [54] E. Jang, S. Gu, and B. Poole, “Categorical reparameterization with gumbel-softmax,” *arXiv preprint arXiv:1611.01144*, 2016.
- [55] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *International conference on machine learning*, 2017, pp. 1126–1135.
- [56] B. Huang, Z. Wang, J. Yang, J. Ai, Q. Zou, Q. Wang, and D. Ye, “Implicit identity driven deepfake face swapping detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 4490–4499.
- [57] D. Afchar, V. Nozick, J. Yamagishi, and I. Echizen, “Mesonet: a compact facial video forgery detection network,” in *2018 IEEE international workshop on information forensics and security (WIFS)*. IEEE, 2018, pp. 1–7.
- [58] H. H. Nguyen, F. Fang, J. Yamagishi, and I. Echizen, “Multi-task learning for detecting and segmenting manipulated facial images and videos,” in *2019 IEEE 10th International Conference on Biometrics Theory, Applications and Systems (BTAS)*. IEEE, 2019, pp. 1–8.
- [59] J. Thies, M. Zollhofer, M. Stamminger, C. Theobalt, and M. Nießner, “Face2face: Real-time face capture and reenactment of rgb videos,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2387–2395.
- [60] <https://github.com/MarekKowalski/FaceSwap/>.
- [61] R. Rothe, R. Timofte, and L. V. Gool, “Deep expectation of real and apparent age from a single image without facial landmarks,” *International Journal of Computer Vision*, vol. 126, no. 2-4, pp. 144–157, 2018.
- [62] J. Deng, J. Guo, E. Ververas, I. Kotsia, and S. Zafeiriou, “Retinaface: Single-shot multi-level face localisation in the wild,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 5203–5212.
- [63] M. Tan and Q. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 6105–6114.
- [64] D. Cozzolino, J. Thies, A. Rössler, C. Riess, M. Nießner, and L. Verdoliva, “Forensicttransfer: Weakly-supervised domain adaptation for forgery detection,” *arXiv preprint arXiv:1812.02510*, 2018.
- [65] H. H. Nguyen, F. Fang, J. Yamagishi, and I. Echizen, “Multi-task learning for detecting and segmenting manipulated facial images and videos,” in *2019 IEEE 10th International Conference on Biometrics Theory, Applications and Systems (BTAS)*. IEEE, 2019, pp. 1–8.
- [66] Y. Qian, G. Yin, L. Sheng, Z. Chen, and J. Shao, “Thinking in frequency: Face forgery detection by mining frequency-aware clues,” in *European Conference on Computer Vision*. Springer, 2020, pp. 86–103.
- [67] D. Li, Y. Yang, Y.-Z. Song, and T. Hospedales, “Learning to generalize: Meta-learning for domain generalization,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, no. 1, 2018.
- [68] K. Sun, H. Liu, Q. Ye, Y. Gao, J. Liu, L. Shao, and R. Ji, “Domain general face forgery detection by learning to weight,” in *Proceedings of the AAAI conference on Artificial Intelligence*, vol. 35, no. 3, 2021, pp. 2638–2646.
- [69] K. Sun, T. Yao, S. Chen, S. Ding, J. Li, and R. Ji, “Dual contrastive learning for general face forgery detection,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 2, 2022, pp. 2316–2324.
- [70] J. Wang, Z. Wu, W. Ouyang, X. Han, J. Chen, Y.-G. Jiang, and S.-N. Li, “M2tr: Multi-modal multi-scale transformers for deepfake detection,” in *Proceedings of the 2022 International Conference on Multimedia Retrieval*, 2022, pp. 615–623.
- [71] M. Tan and Q. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” in *International conference on machine learning*. PMLR, 2019, pp. 6105–6114.
- [72] H. Zhao, W. Zhou, D. Chen, T. Wei, W. Zhang, and N. Yu, “Multi-attentional deepfake detection,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 2185–2194.
- [73] C. Wang and W. Deng, “Representative forgery mining for fake face detection,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 14 923–14 932.
- [74] S. Chen, T. Yao, Y. Chen, S. Ding, J. Li, and R. Ji, “Local relation learning for face forgery detection,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 2, 2021, pp. 1081–1088.
- [75] K. Sun, H. Liu, Q. Ye, Y. Gao, J. Liu, L. Shao, and R. Ji, “Domain general face forgery detection by learning to weight,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, no. 3, 2021, pp. 2638–2646.
- [76] Y. Wang, K. Yu, C. Chen, X. Hu, and S. Peng, “Dynamic graph learning with content-guided spatial-frequency relation reasoning for deepfake detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7278–7287.
- [77] J. Li, D. Li, C. Xiong, and S. Hoi, “Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation,” in *ICML*, 2022.
- [78] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, “Deformable detr: Deformable transformers for end-to-end object detection,” *arXiv preprint arXiv:2010.04159*, 2020.
- [79] S. Woo *et al.*, “Add: Frequency attention and multi-view based knowledge distillation to detect low-quality compressed deepfake images,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 1, 2022, pp. 122–130.