

# Detect Any Deepfakes: Segment Anything Meets Face Forgery Detection and Localization

Yingxin Lai<sup>1</sup>, Zhiming Luo<sup>1</sup> and Zitong Yu<sup>2\*</sup>

<sup>1</sup> Xiamen University

<sup>2</sup> Great Bay University

**Abstract.** The rapid advancements in computer vision have stimulated remarkable progress in face forgery techniques, capturing the dedicated attention of researchers committed to detecting forgeries and precisely localizing manipulated areas. Nonetheless, with limited fine-grained pixel-wise supervision labels, deepfake detection models perform unsatisfactorily on precise forgery detection and localization. To address this challenge, we introduce the well-trained vision segmentation foundation model, i.e., Segment Anything Model (SAM) in face forgery detection and localization. Based on SAM, we propose the Detect Any Deepfakes (DADF) framework with the Multiscale Adapter, which can capture short- and long-range forgery contexts for efficient fine-tuning. Moreover, to better identify forged traces and augment the model’s sensitivity towards forgery regions, Reconstruction Guided Attention (RGA) module is proposed. The proposed framework seamlessly integrates end-to-end forgery localization and detection optimization. Extensive experiments on three benchmark datasets demonstrate the superiority of our approach for both forgery detection and localization. The codes will be released soon at <https://github.com/laiyingxin2/DADF>.

**Keywords:** Deepfake · SAM · Adapter · Reconstruction learning.

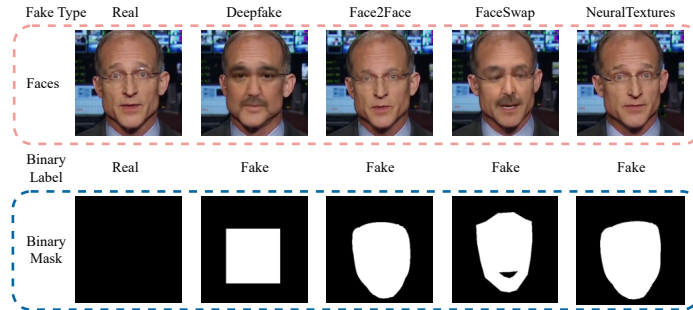
## 1 Introduction

Amongst the diverse human biometric traits, the face is endowed with relatively abundant information and holds significant prominence in identity authentication and recognition. Nonetheless, with the rapid progress of computer vision technology, an array of face-changing techniques has emerged. In particular, the widespread usage of software such as FaceApp and FakeApp has drawn considerable attention to the field of face forgery detection [21]. Therefore, both industry and academia are in urgent need of robust detection methods to mitigate the potential misuse of face forgery technology.

Currently, the majority of forgery detection methods treat the task as a binary classification problem [24,33,7,34] and utilize convolutional neural networks (CNNs) for feature extraction and classification. Although continuous advancements in forged face detection technology in recent years, accurate localization

---

\* Corresponding author



**Fig. 1.** Visualization of the real face and its four common forgery manners in FaceForensics++ [32], as well as their binary labels and corresponding forgery areas.

of forged regions remains a challenge, particularly for models that solely provide classification results. The precise identification of forged regions holds the utmost importance for uncovering the intentions and interpretability of the perpetrators. It allows individuals to discern fake images based on the forged regions and observe the discrepancies between forged and genuine images. Fig. 1 illustrates four common forgery manners in FaceForensics++ [32] and their corresponding pixel-level masks, which are exhausted to be annotated. Due to the limited fine-grained pixel-wise forgery labels, some forgery localization methods [19,17,14,13] trained from scratch usually suffer overfitting.

Recently, Meta introduces the pioneering foundational segmentation model, i.e., Segment Anything Model (SAM) [18,36], demonstrating robust zero-shot segmentation capabilities. Subsequently, researchers have explored diverse approaches such as Low-Rank Adaptation (LoRA) [37], SAM adapter [6], and learnable Prompt [30] to fine-tune SAM on downstream tasks including medical image segmentation and anomaly detection. However, these methods usually yield unsatisfactory positioning outcomes in face forgery localization due to their weak forgery local and global context modeling capacities.

This study focuses on 1) investigating how SAM and its variants perform in the deepfake detection and localization task; and 2) designing accurate and robust pixel-level forgery localization methods across various datasets. For the former one, we find that due to the limited multiscale and subtle forgery context representation capacity, SAM [18] without or with fine-tuning [37,30] cannot achieve satisfactory forgery detection and localization results, which can be alleviated via the proposed SAM based Multiscale Adapter. On the other hand, we find that SAM and its variants are sensitive by the forgery boundary and domain shifts, which might be mitigated by the proposed Reconstruction Guided Attention module. Our main contributions are summarized as follows:

- We are the first to explore the availability of SAM and its fine-tuning strategies in the deepfake detection area. Based on SAM, a novel and efficient Detect Any Deepfakes (DADF) framework is proposed.
- We propose the Multiscale Adapter in SAM, which can capture short- and long-range forgery contexts for efficient fine-tuning.

- We propose the Reconstruction Guided Attention (RGA) module to enhance forged traces and augment the model’s sensitivity towards forgery regions.
- The proposed method achieves state-of-the-art performance in terms of both forgery detection and localization.

## 2 Related work

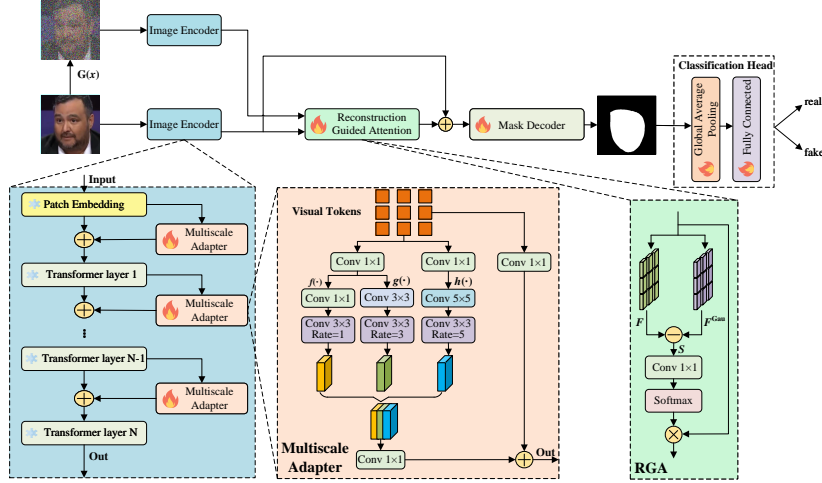
### 2.1 Face Forgery Detection

**Binary classification.** Face forgery detection is predominantly treated as a binary (real/fake) classification task. Plenty of deep learning based methods are developed for detecting face forgery. Li et al. [23] observed distinct variations in blink frequency between forged and authentic videos. To exploit this disparity, the authors utilized CNNs and Long short-term memory (LSTM) [16] models to extract blink-related features in the time domain, enabling the classification of video authenticity. Agarwal et al. [2] observed a mismatch between the movements of the mouth, ears, and chin in fake faces despite synchronized audio. Dang et al [10] incorporated an attention mechanism to emphasize the forgery area, leading to improved accuracy in forgery classification. Alternatively, Nguyen et al. [27] proposed a capsule network designed specifically for identifying counterfeit images or videos. Additionally, Fang et al. [26] introduced the integration of reconstruction losses alongside classification to enhance the overall performance of the model. However, the above-mentioned methods only provide the result of forgery on the scale of the whole image, thus ignoring the identification of the forged region, which lack sufficient interpretability.

**Joint detection and localization.** Face forgery localization precisely identifies manipulated regions of a face at the pixel level. A hybrid CNN-LSTM model [3] was proposed for learning the distinctive boundary variations between manipulated and non-manipulated regions. Nguyen et al. [26] proposed to utilize multi-task learning to detect and locate manipulated regions in both images and videos. Attention mechanisms [10,19] are used to enhance the feature maps of the classification task by highlighting the forged regions. Li et al. [21] proposed to detect the edge regions with mixed boundaries between the manipulated face and the background. However, there are still no works investigating vision segmentation foundation models for joint face forgery detection and localization.

### 2.2 Foundation Model

In recent years, the development of large-scale deep learning pre-training models has promoted the rapid progress of basic visual models in various tasks in the computer field. Models such as BERT [11] and GPT have improved abilities of language understanding, inference and generation in the field of natural language processing, and these models only need the Prompt of specific tasks to apply to new language tasks. CLIP [31] uses contrast loss to learn large-scale image-text pairs. It achieves excellent classification performance in specific downstream tasks without additional data training. COSTA [25] combines the prior knowledge of various pre-training paradigms and unifies them to achieve the most advanced zero-shot classification. DINOv2 [28] can detect objects in



**Fig. 2.** Framework of the proposed Detect Any Deepfakes (DADF). The Multiscale Adapters, Reconstruction Guided Attention module, mask decoder of SAM, and the classification head are trainable while the pre-trained SAM encoder is fixed.

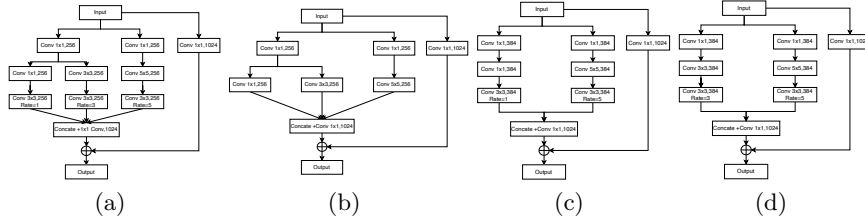
the open world through a given text. SAM [18] proposes a basic model based on image segmentation, using a diverse, high-resolution, licensed and privacy-protected 11 million images and 1.1 billion high-quality segmentation masks for training. The model can accept multiple prompts as input, such as points, boxes and prompts, and can extract high-quality target masks in open-world scenarios. Inspired by these basic models, we try to fine-tune the existing models to achieve the goal of face forgery detection and anomaly localization.

### 2.3 Parameter-Efficient Fine-Tuning

The training cost of adjusting all parameters for foundation models when adapting to downstream tasks may be very high, along with great pressure on storage resources. To alleviate this problem, some researchers have developed efficient fine-tuning strategies. Liu et al. [22] proposed continuous learnable prompts for the frozen language model, which significantly reduces the storage and memory usage required for downstream task adaptation. Clark et al. [9] proposed to train a generative network to sample and replace the original input token, thus reducing the demand for computing resources while ensuring training accuracy. Zhou et al. [37] proposed the Low Rank Adaptation (LoRA) to indirectly train few dense layers in transformer layers by optimizing the rank decomposition matrix of the dense layer changes during the adaptation process, while maintaining the weight of the pre-training unchanged. The most similar works to ours are the adapter-based fine-tuning [6,35]. Instead of exploiting only channel-wise new knowledge for SAM, the proposed multiscale adapter is able to mine short- and long-range forgery contexts for efficient SAM fine-tuning.

## 3 Methodology

As illustrated in Fig. 2, the proposed SAM-based [18] architecture involves an image encoder with the Multiscale Adapters for feature extraction, a Reconstruction Guided Attention (RGA) module for forged feature refinement and a



**Fig. 3.** (a) The original Multiscale Adapter; (b) Without dilated Convolutions; (c) Without  $3 \times 3$  convolution branch; (d) Without  $1 \times 1$  Convolution Branch. A batch normalization layer and a ReLU layer are cascaded after every convolution operator. Conv1  $\times$  1,256 means using  $1 \times 1$  convolution with 256 channels.

mask decoder for forgery mask prediction. Based on the predicted mask, a classification head consisting of global average pooling and fully connected layers is cascaded for real/fake classification.

### 3.1 Multiscale Adapter for SAM

In consideration of the limited data scale in the face forgery detection task, we freeze the main parameters (i.e., Transformer layers) of the SAM encoder and insert learnable specific modules to mine task-aware forgery clues. Specifically, we incorporate a concise and efficient Multiscale Adapter module along each transformer layer to capture forgery clues with diverse receptive fields using a multi-scale fashion.

First, we tokenize the input image  $x$  into visual tokens  $x_p = P(x)$  via the Patch Embedding layer (denoted as  $P(\cdot)$ ). During this process, defaulted patch size  $14 \times 14$  is used. Then  $N$  fixed Transformer layers  $L_i, i \in \{1, \dots, N\}$  with learnable Multiscale Adapter  $R(\cdot)$  are used for extracting short- and long-range contextual features  $Z_{\text{tran}}$ , which is then mapped to task-aware forgery features  $F$  via a dimension-matched learnable linear Task Head  $T(\cdot)$ . The feature encoding procedure can be formulated as

$$\begin{aligned} Z_{\text{tran}} &= L_n(\dots L_2(R(L_1(R(x_p))))), \\ F &= T(Z_{\text{tran}}). \end{aligned} \quad (1)$$

As for the Multiscale Adapter (see the middle red block in Fig. 2 and Fig. 5(a)), the output features  $x'$  of each Transformer layer are passed by a  $1 \times 1$  convolution, and then split into three branches  $f(\cdot)$ ,  $g(\cdot)$ , and  $h(\cdot)$ . Each branch uses different convolution kernel sizes and dilated rates for complementary forgery context mining. Therefore, the multiscale short- and long-range features  $Sout_1$  can be formulated as

$$\begin{aligned} Sout' &= \text{Concat}(f(\text{Conv}_{1 \times 1}(x')), g(\text{Conv}_{1 \times 1}(x')), h(\text{Conv}'_{1 \times 1}(x'))), \\ Sout_1 &= \text{Conv}''_{1 \times 1}(Sout'), \end{aligned} \quad (2)$$

where  $f(\cdot)$ ,  $g(\cdot)$  and  $h(\cdot)$  denote a convolution with kernel size  $1 \times 1$  cascaded with a convolution with kernel size  $3 \times 3$  and dilated rates 1, a convolution with kernel size  $3 \times 3$  cascaded with a convolution with kernel size  $3 \times 3$  and dilated rates 3, and a convolution with kernel size  $5 \times 5$  cascaded with a convolution with kernel size  $3 \times 3$  and dilated rates 5, respectively.  $Sout'$  is the result of merging features

of the three branches, which is then re-projected to the original channel size via a  $1 \times 1$  convolution to obtain  $Sout_1$ .

Finally, the resultant multi-scale features  $Sout_1$  are added together to the original features  $x'$  passed over a  $1 \times 1$  convolution operation  $Conv'''_{1 \times 1}$ , ensuring the preservation of the original information. The final multiscale contextual features  $Sout$  can be formulated as

$$Sout = Sout_1 + Conv'''_{1 \times 1}(x'). \quad (3)$$

In terms of the structure of Multiscale Adapter, as shown in Fig. 3, three kinds of variants (i.e., Multiscale Adapter-B, Multiscale Adapter-C, Multiscale Adapter-D) are also investigated according to different scalabilities of receptive fields. The ablation study among them can be found in Table 5.

### 3.2 Reconstruction Guided Attention

In order to enhance the sensitivity to deep forged regions and explore the common and compact feature patterns of real faces, we propose a reconstruction learning method, namely Reconstruction Guided Attention (RGA).

In the training process, we simulate the forged faces by introducing white noise  $G(\cdot)$  on the real faces. Based on the noisy inputs, the model gradually performs feature reconstruction to obtain the reconstructed features  $F^{Gau}$ .

$$\begin{aligned} x^{Gau} &= G(x), \\ F^{Gau} &= \Phi(x^{Gau}), \end{aligned} \quad (4)$$

where  $\Phi$  denotes the whole image encoder. After the feature reconstruction process performed by the image encoder, we compute the absolute difference  $S$  between the original features and the reconstructed features. In this calculation, the function  $|\cdot|$  represents the absolute value function.

$$S = |F^{Gau} - F|. \quad (5)$$

Subsequently, an enhancer  $\varphi(\cdot)$  with  $1 \times 1$  convolution is employed to highlight and enhance regions that might contain forgeries, which is cascaded with the Softmax function layer  $\alpha(\cdot)$  to generate the forgery-aware attention map. Finally, we perform element-wise multiplication based refinement operation  $\otimes$  between the obtained attention weights and the original features to obtain the final features  $F_{\text{final}}$ , which are then sent for the mask decoder. The procedure can be formulated as:

$$F_{\text{final}} = [\alpha(\varphi(S)) \otimes \varphi(F)] + F. \quad (6)$$

After obtaining the features  $F$  from the real faces and the reconstructed features  $F^{Gau}$  from the anomalies/forgeries simulation, we calculate the reconstruction loss  $\mathcal{L}_{\text{rec}}$  for each batch  $M$  with  $L1$  norm. Notably, the reconstruction

loss  $\mathcal{L}_{\text{rec}}$  is exclusively trained on real samples. Ablation studies on calculating  $\mathcal{L}_{\text{rec}}$  for fake faces and all (real+fake) faces can be found in Table 6.

$$\mathcal{L}_{\text{rec}} = \frac{1}{M} \sum_{i \in \mathcal{M}} \|F_i^{\text{Gau}} - F_i\|_1. \quad (7)$$

In the training stage, the RGA module leverages the abnormal simulated faces as one of the inputs and gradually recovers the intrinsic features of the real faces. Through this reconstruction process, SAM models can better understand the common and compact feature patterns of real faces, and even pay more attention to unknown forged regions in the inference stage.

### 3.3 Loss Function

The overall loss function  $\mathcal{L}_{\text{overall}}$  of DADF consists of three components: segmentation loss  $\mathcal{L}_{\text{seg}}$ , classification loss  $\mathcal{L}_{\text{cls}}$ , and feature reconstruction loss  $\mathcal{L}_{\text{rec}}$ . The segmentation loss  $\mathcal{L}_{\text{seg}}$  represents the semantic loss, while the binary cross-entropy loss  $\mathcal{L}_{\text{cls}}$  measures the binary real/fake classification error. The feature reconstruction loss  $\mathcal{L}_{\text{rec}}$  captures the reconstruction error.

$$\mathcal{L}_{\text{overall}} = \mathcal{L}_{\text{seg}} + \lambda_1 \mathcal{L}_{\text{rec}} + \lambda_2 \mathcal{L}_{\text{cls}}, \quad (8)$$

where the hyperparameters  $\lambda_1$  and  $\lambda_2$  are used to balance the different components of the loss, which are set to 0.1 according to empirical observations.

## 4 Experiments

### 4.1 Datasets and Performance Metrics

**Datasets.** The **FaceForensics++** (FF++) [32] utilizes four different algorithms: Deepfakes (DF), Face2Face (F2F), FaceSwap (FS) and NeuralTextures (NT) to generate forgery faces. The video data also provide versions with different compression ratios: original quality (quantization = 0), high quality (HQ, quantization = 23) and low quality (LQ, quantization = 40). The **DF-TIMI** [20] dataset contained 16 pairs of similar people, each of whom lived 10 videos. It includes 960 videos: 320 genuine and 640 forged using FaceSwap technology. Among the forged videos, 320 are high quality (HQ) and 320 are low quality (LQ). The **DFD** [32] was created specifically for DeepFake technology. These videos were procured from the YouTube platform, consisting of 363 authentic videos and 3068 fabricated videos, which have been further categorized as high-quality (HQ) and low-quality (LQ). The **FMLD** [19] comprises 40,000 synthetic faces generated using StyleGAN and PGGAN models, and an additional 40,000 attribute-manipulated faces using StarGAN and AttGAN. Specifically, attribute manipulations include glasses and hair regions.

**Evaluation metrics.** Two commonly used metrics, namely Binary Classification Accuracy (BCA) and Inverse Intersection Non-Containment (IINC) [10] are employed for forgery localization. For fair comparisons, we follow the same evaluation protocols as [19] for face forgery localization. In terms of evaluating the performance of face forgery detection, Accuracy (ACC) is adopted.

Table 1: Results of the forgery localization on FF++ (HQ) [32] and FMLD [19].

Dataset	FF++ (HQ)		FMLD	
	PBCA(%) $\uparrow$	IINC(%) $\downarrow$	PBCA(%) $\uparrow$	IINC(%) $\downarrow$
Multitask [26]	94.88	4.46	98.59	3.52
DFFD Reg [10]	94.85	4.57	98.72	3.31
DFFD Mam [10]	91.45	13.09	96.86	23.93
Locate [19]	95.77	3.62	99.06	<b>2.53</b>
SAM [18]	92.97	4.25	97.29	3.40
SAM+LoRA [37]	93.12	4.78	98.06	3.51
SAM+Prompt [30]	94.60	3.48	98.01	2.82
<b>DADF (Ours)</b>	<b>96.64</b>	<b>3.21</b>	<b>99.26</b>	2.64

Table 2: The forgery detection performance (ACC(%)) on FF++ (LQ) [32].

Methods	DF	FF	FS	NT	Average
Steg. Features [12]	67.00	48.00	49.00	56.00	55.00
Cozzolino [5]	75.00	56.00	51.00	62.00	61.00
Bayar & Stamm [4]	87.00	82.00	74.00	74.00	79.25
Rahmouni [1]	80.00	62.00	59.00	59.00	65.00
MesoNet [22]	90.00	83.00	83.00	75.00	82.75
SPSL [8]	93.48	86.02	92.26	<b>92.26</b>	91.00
Xception [8]	97.16	91.02	96.71	82.88	91.94
Locate [19]	97.25	94.46	97.13	84.63	93.36
SAM [18]	89.32	84.56	91.19	80.01	86.27
SAM+LoRA [37]	90.12	85.41	91.28	80.15	86.74
SAM+Prompt [30]	97.34	95.84	97.44	84.72	93.83
<b>DADF (Ours)</b>	<b>99.02</b>	<b>98.92</b>	<b>98.23</b>	87.61	<b>95.94</b>

## 4.2 Implementation Details

We use the SAM-based [18] ViT-H model as the backbone with a null input prompt setting. We train models with a batch size of 4 and adopt the AdamW optimizer. We employ the cosine decay method. The initial learning rate is set to 0.05. For the FF++ dataset, we conduct 30 epochs of training, while the FMLD dataset requires 50 epochs to train the model effectively. To adapt the learning rate, a step learning rate scheduler is employed. As for the RGA module, white noise is employed as a noise source, which is incorporated into the data using a normal distribution with a zero mean and a variance of 1e-6.

## 4.3 Intra-dataset Testing

**Results of face forgery localization.** Table 1 presents the results of the forgery localization on FF++ (HQ) [32] and FMLD [19]. The proposed DADF outperform the classical face forgery localization model [19] by 0.87% and 0.2% PBCA on FF++ (HQ) and FMLD, respectively. We can also find from the results of SAM [18] that direct finetuning SAM cannot achieve acceptable face forgery localization performance due to its heavy model parameters and limited task-aware data. Despite slight improvement via parameter-efficient fine-tuning strategies, SAM with LoRA or Prompt still has performance gaps with the previous localization method [19]. Thanks to the rich forgery contexts from the Multiscale Adapter and the strong forgery attention ability of RGA module, the proposed DADF improves baseline SAM [18] by 3.67%/-1.04% and 1.97%/-0.76% PBCA/IINC on FF++ (HQ) and FMLD, respectively.

**Results of face forgery detection.** Table 2 presents the detection accuracy (ACC) of our model on various forgery techniques, namely Deepfake (DF), Face2Face (FF), FaceSwap (FS), and NeuralTextures (NT), using the challenging FF++ (LQ) [32]. It is clear that the proposed DADF performs significant



Table 3: Results of cross-dataset face forgery detection.

Dataset	DFD (LQ)		DF-TIMIT (HQ)		DF-TIMIT (LQ)		Average	
	AUC(%) $\uparrow$	EER(%) $\downarrow$	AUC(%) $\uparrow$	EER(%) $\downarrow$	AUC(%) $\uparrow$	EER(%) $\downarrow$	AUC(%) $\uparrow$	EER(%) $\downarrow$
Method	52.25	48.65	33.61	60.16	45.08	53.04	34.64	53.95
MesoNet [1]								
MesoIncep4 [1]	<b>63.27</b>	40.37	16.12	76.18	27.47	66.77	35.62	61.10
ResNet50 [15]	60.61	42.23	41.95	55.97	47.27	52.33	49.94	50.17
Face X-ray [21]	62.89	39.58	42.52	55.07	50.05	<b>49.11</b>	51.81	47.92
DFFD [10]	60.60	42.32	32.91	61.16	39.32	57.06	44.27	53.51
Multi-task [26]	58.61	44.49	16.53	77.86	15.59	78.50	30.24	66.95
F3Net [29]	58.89	39.87	29.12	58.33	45.67	52.72	44.56	50.30
Xception [8]	59.73	43.12	33.82	62.83	40.79	57.44	44.78	54.46
SAM [18]	50.61	49.13	43.19	57.94	45.71	54.39	46.50	53.82
SAM+LoRA [37]	53.71	48.29	43.64	56.67	47.64	53.02	48.33	52.66
SAM+Prompt [30]	57.25	45.28	44.32	55.07	48.17	52.54	49.91	50.96
<b>DADF (Ours)</b>	63.21	<b>39.52</b>	<b>46.37</b>	<b>53.20</b>	<b>50.62</b>	49.74	<b>53.40</b>	<b>47.48</b>

Table 4: Ablation studies on the FF++ (HQ) [32] dataset.

Baseline (SAM)	Multiscale Adapter	RGA	Localization		Detection
			PBLA(%) $\uparrow$	IINC(%) $\downarrow$	ACC(%) $\uparrow$
$\checkmark$			92.97	4.25	86.27
$\checkmark$	$\checkmark$		96.31	3.40	94.48
$\checkmark$		$\checkmark$	95.61	3.96	92.64
$\checkmark$	$\checkmark$	$\checkmark$	<b>96.64</b>	<b>3.21</b>	<b>95.94</b>

improvements in classification accuracy compared to previous methods among different forgery techniques. This highlights the effectiveness of our Multiscale Adapter and RGA module in enhancing the detection capabilities, compared with the original SAM [18] and its variants (SAM+LoRA [37] and SAM+Prompt [30]). Specifically, the proposed DADF improves more than 3% ACC compared with the second-best method on Face2Face detection.

#### 4.4 Cross-dataset Testing

In order to assess the generalization ability of our method on unseen domains and unknown deepfakes, we conducted cross-dataset experiments by training and testing on different datasets. Specifically, we train models on FF++ (LQ), and then test them on DFD (LQ), DF-TIMIT (HQ), and DF-TIMIT (LQ). The results shown in Table 3 demonstrate that the proposed DADF outperforms other methods in terms of average performance among the three testing settings. Compared with SAM and its variants, the significant improvement of DADF in performance is attributed to 1) the introduction of the Multiscale Adapter, which enables forgery feature learning from diverse receptive fields; and 2) the attentional forgery feature refinement via the RGA module, enhancing the robustness under domain shifts and perception of forgery regions.

#### 4.5 Ablation Study

To validate the effectiveness of the Multiscale Adapter and Reconstruction Guided Attention module, ablation experiments are conducted on FF++ (HQ).

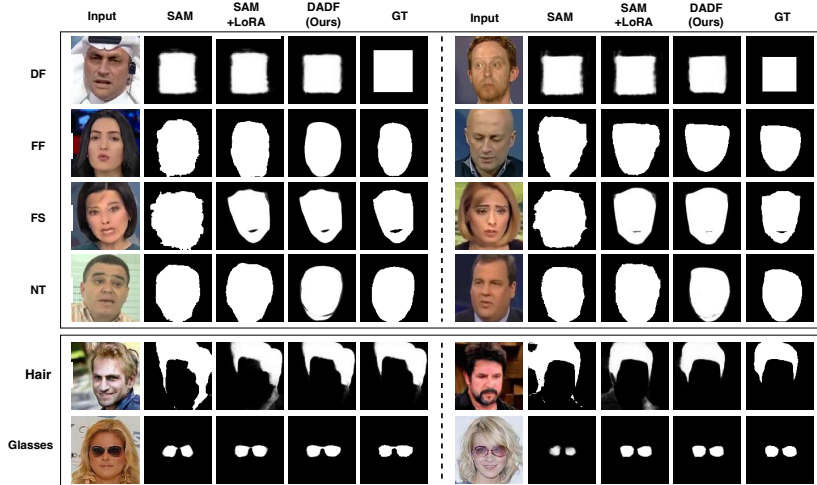
**Efficacy of the Multiscale Adapter.** It can be seen from the first two rows of Table 4 that compared with baseline SAM-only fine-tuning, SAM with Multiscale

Table 5: Ablations of different Multiscale Adapters on FF++ (HQ) [32] dataset.

Module	Localization		Detection
	PBLA(%) $\uparrow$	IINC(%) $\downarrow$	ACC(%) $\uparrow$
Multiscale Adapter-B	96.52	3.43	95.79
Multiscale Adapter-C	96.57	3.31	95.87
Multiscale Adapter-D	96.61	3.26	95.91
Multiscale Adapter	<b>96.64</b>	<b>3.21</b>	<b>95.94</b>

Table 6: Ablation studies of  $\mathcal{L}_{\text{rec}}$  calculation on FF++ (HQ) [32] dataset.

Data	Localization		Detection
	PBLA(%) $\uparrow$	IINC(%) $\downarrow$	ACC(%) $\uparrow$
Real & Fake	92.15	3.62	91.42
Fake	93.34	3.56	92.16
Real	<b>96.64</b>	<b>3.21</b>	<b>95.94</b>



**Fig. 4.** Visualization of face forgery localization results of various methods on the FF++ (HQ) dataset (DF, FF, FS, and NT) [32] and FMLD (Hair and Glasses) [19].

Adapter improves 3.34%/-0.85% PBLA/IINC for forgery localization and 8.21% ACC for forgery detection on the FF++ (HQ). Considering different configurations (see Fig. 3) of Multiscale Adapter, Table 5 demonstrates that removing dilated convolutions results in the largest accuracy drop, while the best performance is achieved by incorporating multiscale convolution modules alongside dilated convolutions.

**Efficacy of the RGA.** As shown in the last two rows of Table 4, equipping with RGA module can improve the baseline SAM by 2.64%/-0.29% PBLA/IINC for forgery localization and 6.37% ACC for forgery detection on the FF++ (HQ). Similarly, based on the SAM with Multiscale Adapter, the RGA module can further benefit the forgery localization by 0.33% PBLA and detection by 1.46% ACC. As for the loss function  $\mathcal{L}_{\text{rec}}$  calculation for RGA, it can be seen from Table 6 that the performance drops sharply when  $\mathcal{L}_{\text{rec}}$  calculated for fake faces and all (real+fake) faces, which might result from the redundant features of real faces and less attention on anomalies.

#### 4.6 Visualization and Discussion

We visualize some representative forgery samples and their mask labels as well as predictions in Fig. 4. It is evident that the forgery localization quality from the proposed DADF outperforms SAM and its LoRA fine-tuning in accurately localizing and closely resembling the ground truth, particularly in fine-grained details such as edge, boundary, and face-head contexts.

Besides, the proposed Multiscale Adapter is a parameter-efficient fine-tuning strategy alternative to tune the entire Transformer layers. Remarkably, by adjusting only 18.64% parameters of the SAM, substantial benefits on face forgery detection and localization are achieved, including reduced training costs and improved practical performance.

## 5 Conclusion

In this paper, we introduce a Segment Anything Model based face forgery detection and localization framework, namely Detect Any Deepfakes (DADF). Specifically, we propose the Multiscale Adapter and Reconstruction Guided Attention (RGA) to efficiently fine-tune SAM with rich contextual forgery clues and enhance the robustness of forgery localization. Extensive experimental results validate the effectiveness of the proposed DADF across different qualities of face images and even under cross-domain scenarios.

## References

1. Afchar, D., Nozick, V., Yamagishi, J., Echizen, I.: Mesonet: a compact facial video forgery detection network. In: IEEE WIFS (2018)
2. Agarwal, S., Farid, H.: Detecting deep-fake videos from aural and oral dynamics. In: IEEE CVPR (2021)
3. Bappy, J.H., Roy-Chowdhury, A.K., Bunk, J., Nataraj, L., Manjunath, B.: Exploiting spatial structure for localizing manipulated image regions. In: IEEE ICCV (2017)
4. Bayar, B., Stamm, M.C.: A deep learning approach to universal image manipulation detection using a new convolutional layer. In: ACM workshop on information hiding and multimedia security (2016)
5. Cao, J., Ma, C., Yao, T., Chen, S., Ding, S., Yang, X.: End-to-end reconstruction-classification learning for face forgery detection. In: CVPR (2022)
6. Chen, T., Zhu, L., Ding, C., Cao, R., Zhang, S., Wang, Y., Li, Z., Sun, L., Mao, P., Zang, Y.: Sam fails to segment anything? – sam-adapter: Adapting sam in underperformed scenes: Camouflage, shadow, and more (2023)
7. Chierchia, G., Parrilli, S., Poggi, G., Verdoliva, L., Sansone, C.: Prnu-based detection of small-size image forgeries. In: IEEE DSP (2011)
8. Chollet, F.: Xception: Deep learning with depthwise separable convolutions. In: IEEE CVPR (2017)
9. Clark, K., Luong, M.T., Le, Q.V., Manning, C.D.: Electra: Pre-training text encoders as discriminators rather than generators. arXiv preprint arXiv:2003.10555 (2020)
10. Dang, H., Liu, F., Stehouwer, J., Liu, X., Jain, A.K.: On the detection of digital face manipulation. In: IEEE CVPR (2020)
11. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
12. Fridrich, J., Kodovsky, J.: Rich models for steganalysis of digital images. IEEE TIFS (2012)

13. Guan, W., Wang, W., Dong, J., Peng, B., Tan, T.: Collaborative feature learning for fine-grained facial forgery detection and segmentation. arXiv preprint arXiv:2304.08078 (2023)
14. Guo, X., Liu, X., Ren, Z., Grosz, S., Masi, I., Liu, X.: Hierarchical fine-grained image forgery detection and localization. In: IEEE CVPR (2023)
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: IEEE CVPR (2016)
16. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural computation (1997)
17. Huang, Y., Juefei-Xu, F., Guo, Q., Liu, Y., Pu, G.: Fakelocator: Robust localization of gan-based face manipulations. IEEE TIFS (2022)
18. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al.: Segment anything. arXiv preprint arXiv:2304.02643 (2023)
19. Kong, C., Chen, B., Li, H., Wang, S., Rocha, A., Kwong, S.: Detect and locate: Exposing face manipulation by semantic-and noise-level telltales. IEEE TIFS (2022)
20. Korshunov, P., Marcel, S.: Deepfakes: a new threat to face recognition? assessment and detection. arXiv preprint arXiv:1812.08685 (2018)
21. Li, L., Bao, J., Zhang, T., Yang, H., Chen, D., Wen, F., Guo, B.: Face x-ray for more general face forgery detection. In: IEEE CVPR (2020)
22. Liu, X., Ji, K., Fu, Y., Tam, W.L., Du, Z., Yang, Z., Tang, J.: P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. arXiv preprint arXiv:2110.07602 (2021)
23. LIY, C.M., InIctuOculi, L.: Exposingaicreated fakevideosbydetectingeyebinking. In: IEEE WIFS (2018)
24. Lukáš, J., Fridrich, J., Goljan, M.: Detecting digital image forgeries using sensor pattern noise. In: Security, steganography, and watermarking of multimedia contents VIII. SPIE (2006)
25. Mensink, T., Gavves, E., Snoek, C.G.: Costa: Co-occurrence statistics for zero-shot classification. In: CVPR (2014)
26. Nguyen, H.H., Fang, F., Yamagishi, J., Echizen, I.: Multi-task learning for detecting and segmenting manipulated facial images and videos. In: IEEE BTAS (2019)
27. Nguyen, H.H., Yamagishi, J., Echizen, I.: Capsule-forensics: Using capsule networks to detect forged images and videos. In: IEEE ICASSP (2019)
28. Oquab, M., Darcet, T., Moutakanni, T., Vo, H., Szafraniec, M., Khalidov, V., Fernandez, P., Haziza, D., Massa, F., El-Nouby, A., et al.: Dinov2: Learning robust visual features without supervision. arXiv preprint arXiv:2304.07193 (2023)
29. Qian, Y., Yin, G., Sheng, L., Chen, Z., Shao, J.: Thinking in frequency: Face forgery detection by mining frequency-aware clues. In: ECCV (2020)
30. Qiu, Z., Hu, Y., Li, H., Liu, J.: Learnable ophthalmology sam. arXiv preprint arXiv:2304.13425 (2023)
31. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: ICML (2021)
32. Rössler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., Nießner, M.: Face-forensics: A large-scale video dataset for forgery detection in human faces. arXiv preprint arXiv:1803.09179 (2018)
33. Swaminathan, A., Wu, M., Liu, K.R.: Digital image forensics via intrinsic fingerprints. IEEE TIFS (2008)
34. Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: ICML (2019)

35. Wu, J., Fu, R., Fang, H., Liu, Y., Wang, Z., Xu, Y., Jin, Y., Arbel, T.: Medical sam adapter: Adapting segment anything model for medical image segmentation. arXiv preprint arXiv:2304.12620 (2023)
36. Zhang, C., Liu, L., Cui, Y., Huang, G., Lin, W., Yang, Y., Hu, Y.: A comprehensive survey on segment anything model for vision and beyond. arXiv preprint arXiv:2305.08196 (2023)
37. Zhou, X., Yang, C., Zhao, H., Yu, W.: Low-rank modeling and its applications in image analysis. ACM Computing Surveys (CSUR) (2014)