# Enhancing General Face Forgery Detection via Vision Transformer with Low-Rank Adaptation

Chenqi Kong, Haoliang Li, and Shiqi Wang

## Abstract

*Nowadays, forgery faces pose pressing security concerns over fake news, fraud, impersonation, etc. Despite the demonstrated success in intra-domain face forgery detection, existing detection methods lack generalization capability and tend to suffer from dramatic performance drops when deployed to unforeseen domains. To mitigate this issue, this paper designs a more general fake face detection model based on the vision transformer(ViT) architecture. In the training phase, the pretrained ViT weights are freezed, and only the Low-Rank Adaptation(LoRA) modules are updated. Additionally, the Single Center Loss(SCL) is applied to supervise the training process, further improving the generalization capability of the model. The proposed method achieves state-of-the-arts detection performances in both cross-manipulation and cross-dataset evaluations.*

## 1. Introduction

With the rapid proliferation of digital face medias circulating on social media, non-expert attackers can easily create fake face content due to unrestricted access to face media and the ease of implementing face manipulation techniques (*e.g.*, Deepfakes [2], Face2Face [3], FaceSwap [40], NeuralTextures [39], and other attacks [43, 42].) [20]. Even worse, the availability of commercial tools and products (e.g., FakeApp [1]) makes generating forgery faces much easier. The abuse of face manipulation have posed grand security concerns to the public at large, including fake news, financial fraud, identity theft, etc [19]. Thus, it is of utmost importance to propose detection methods to counter the malicious attacks[6, 7] and build trust of digital facial medias.

The past decades have witnessed significant progress in face forgery detection methodologies. Early works mainly focus on extracting handcrafted features such as lack of eye blinking [27], head pose inconsistency [41], and face warping artifacts [28] from the inputs. However, these methods suffer from limited accuracy and low generalization capability. Thanks to the advent of artificial intelligence and deep learning, many learning-based detection methods [34, 32, 10, 26, 18, 30] have been proposed and achieved outstanding detection performance under intra-domain settings. Nonetheless, learning-based methods are prone to overfitting to the training data, resulting in dramatic performance drops when deployed to unforeseen domains. In this vein, follow-up works such as [25, 37, 36, 31] aim to mine more inherent and general artifacts from different manipulation techniques and datasets. Some multimodal-based models seek to use auxiliary modalities (*e.g.*, audio modality) for more robust defense [44, 22, 13, 21].

Inspired by the recent success of vision transformer (ViT) [12], we apply the powerful ViT as the backbone of our framework. To achieve more general face forgery detection performance, we propose to incorporate the Low-Rank Adaptation (LoRA) [16] in this method. LoRA is a parameter-efficient tuning method that has been demonstrated effective in various domain generalization and few-shot learning tasks. Moreover, we borrow the idea of Single-Center Loss(SCL) [25] to make the features of real faces more compact and push the fake features away from the center of real features, thereby achieving more general face forgery detection. Compared with the ViT baseline, the designed model achieves a 6.6% and 11.19% AUC score boosts in challenging low-quality cross-manipulation and cross-dataset evaluations, respectively.

## 2. Related Work

In this section, we first provide a broad review of prior literature on face forgery detection. Then, we briefly analyze and discuss typical parameter-efficient ViT tuning methods.

### 2.1 Face forgery detection methods

Early works on forgery face detection focused on extracting handcrafted features from the input face images/videos. Li *et.al.* [27] analyzed the eye-blinking frequency to identify input authentication, while follow-up works [41, 28] detected the head-pose inconsistency and face warping artifacts to determine the input face videos as real *v.s.* fake. With the advent of artificial intelligence and deep learning, numerous learning-based methods have been

proposed and achieved promising detection accuracy. Qian *et.al.* [34] propose to mine forgery-related features in frequency domain and achieved promising classification accuracy for low-quality fake videos. Dang *et.al.* [10] proposed using manipulation region maps and applying an attention mechanism to achieve more accurate performance. Kong *et.al.* [18] exploited both manipulation region and noise map to supervise the model training and obtain outstanding detection performance. Despite their demonstrated success in intra-domain evaluations, most existing face forgery detection methods lack generalization capability and cannot adapt well when deployed in unseen environments. To mitigate this issue, Face X-ray [26] proposed highlighting the boundary of manipulation regions in fake faces, thus achieving more general detection performance over different forgery techniques. SBI [36] proposed a novel data augmentation method only using real face images to enforce the model to focus on inherent forgery artifacts rather than semantic contents. LTW [37] designed a more general model using meta learning and achieves outstanding cross-domain detection performance. This paper aims to achieve more general face forgery detection by taking advantage of LoRA modules and single-center loss to realize accurate and robust face manipulation detection.

## 2.2 Parameter-efficient tuning for ViT

In the past two years, vision transformers(ViT) have seen great success and have exploded into a plethora of vision applications, such as image classification, semantic segmentation, and object detection. Pretrained ViT models have also been widely used in downstream tasks and have achieved outstanding results through transfer learning. To improve the generalization capability of ViT and reduce computational cost, several parameter-efficient tuning methodologies have been proposed, such as Adapter [15], Low-Rank Adaptation (LoRA) [16], and Visual Prompt Tuning (VPT) [17]. ViT Adapter is a neural network that includes a down-sample and up-sample layer, while LoRA optimizes the rank-decomposed changes of the two projection layers. VPT can be regarded as extra learnable input tokens in input space. Typically, these modules will be tuned and the ViT backbone parameters will be frozen with pretrained weights in the training phase. While the parameter-efficient tuning modules have been widely used in domain generalization and few-shot learning, how the tuning modules benefit the general forgery face detection has not been investigated yet. In this paper, we investigate the effectiveness of LoRA in forgery face detection and suggest incorporating Adapter and VPT in future works.
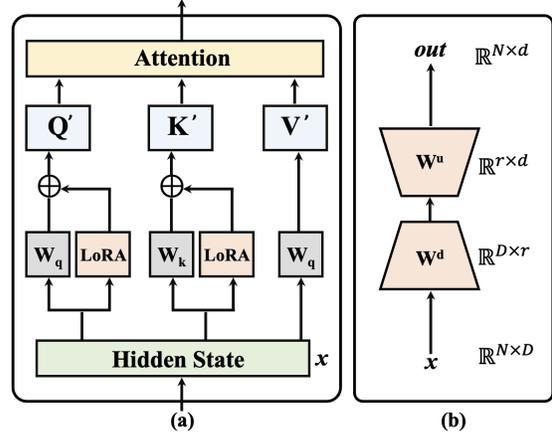


Figure 1: (a). Illustration of the attention mechanism assembled with LoRA; (b). Details of LoRA.

## 3. Proposed Method

In this paper, we take the ViT [12] as our backbone. As shown in Fig. 1, we apply LoRA to the weights of query and key, in each attention layer. We fix the ViT with ImageNet weights and only update the LoRA parameters during the training process. The objective function of the model is the weighted summation of the cross-entropy loss $L_{ce}$ and the single-center loss $L_{scl}$:

$$L = L_{ce} + \lambda L_{scl}, \tag{1}$$

where $\lambda$ is the loss weight. The cross-entropy loss is detailed as:

$$L_{ce} = -\frac{1}{N} \sum_{i=1}^{N} (c_i \log \hat{c}_i + (1 - c_i) \log(1 - \hat{c}_i)), \tag{2}$$

where N is the number of input images. $\hat{c}_i$ and $c_i$ are the prediction result and ground-truth label. We dedicate the single-center loss $L_{scl}$ in Sec. 3.2.

### 3.1. Low-Rank Adaptation(LoRA)

In typical attention mechanism, the query Q, key K, and value V can be obtained via Eqn.(3):

$$Q = W_q x, K = W_k x, V = W_v x \tag{3}$$

where $x \in \mathbb{R}^{N \times D}$, $(W_q, W_k, W_v) \in \mathbb{R}^{D \times d}$. $W_q$, $W_k$, and $W_v$ are learnable weights.

In this model, we tune Low-Rank Adaptation(LoRA) modules, instead of $W_q$, $W_k$, and $W_v$, to obtain more general results. Fig. 1 (a) illustrates the modified attention mechanism. $Q'$, $K'$, and $V'$ can be obtained by Eqn.(4):

$$Q' = W_q x + s W_q^d W_q^u x; K' = W_k x + s W_k^d W_k^u x; V' = W_v x \tag{4}$$

Table 1: High quality (c23) cross-manipulation detection performance on unseen forgery methods.

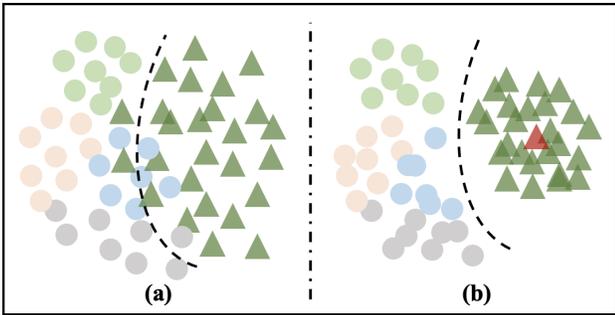| Setting | FF,FS,NT→DF | | DF,FS,NT→FF | | DF,FF,NT→FS | | DF,FF,FS→NT | | Average | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Method | AUC | ACC | AUC | ACC | AUC | ACC | AUC | ACC | AUC | ACC |
| ResNet18 [14] | 0.813 | 0.656 | 0.746 | 0.596 | 0.464 | 0.476 | 0.688 | 0.528 | 0.678 | 0.564 |
| Xception [8] | 0.907 | 0.795 | 0.753 | 0.558 | 0.460 | 0.472 | 0.744 | 0.557 | 0.716 | 0.596 |
| EfficientNet [38] | 0.485 | 0.495 | 0.556 | 0.523 | 0.517 | 0.517 | 0.493 | 0.500 | 0.513 | 0.509 |
| All-train EfficientNet [38] | 0.911 | 0.824 | 0.801 | 0.633 | 0.543 | 0.500 | 0.774 | 0.608 | 0.757 | 0.641 |
| Focal-loss EfficientNet [38] | 0.903 | 0.813 | 0.798 | 0.608 | 0.503 | 0.484 | 0.759 | 0.604 | 0.741 | 0.627 |
| Forensics Transfer [9] | N.A. | 0.720 | N.A. | 0.645 | N.A. | 0.460 | N.A. | 0.569 | N.A. | 0.599 |
| Multi-task [33] | N.A. | 0.703 | N.A. | 0.587 | N.A. | 0.497 | N.A. | 0.603 | N.A. | 0.598 |
| MLDG [24] | 0.918 | 0.842 | 0.771 | 0.634 | 0.609 | 0.527 | **0.780** | 0.621 | 0.770 | 0.656 |
| LTW [37] | 0.927 | 0.856 | 0.802 | 0.656 | 0.640 | 0.549 | 0.773 | **0.653** | 0.786 | 0.679 |
| ViT base [12] | 0.771 | 0.701 | 0.656 | 0.582 | 0.510 | 0.498 | 0.554 | 0.517 | 0.623 | 0.575 |
| Ours | **0.935** | **0.862** | **0.875** | **0.753** | **0.651** | **0.554** | 0.707 | 0.626 | **0.792** | **0.699** |



**(a)** | **(b)**

Figure 2: Illustration of feature distribution (a). without single-center loss; (b). with single-center loss.

where $s$ is the fixed scale parameter. Fig. 1 (b) shows the LoRA details. $W_d \in \mathbb{R}^{D \times r}$, $W_u \in \mathbb{R}^{r \times d}$, and $r$ is a hyper-parameter and generally much smaller than $d$ and $D$. By optimizing their rank-decomposed changes $W^d W^u$, we can benefit from the LoRA modules in following two aspects: (a). the proposed architecture is more computational-efficient since the number of the trainable parameters has been greatly reduced (going from $3Dd$ to $r(D + d)$); (b). the model retains abundant knowledge learned from ImageNet dataset and can be flexibly transferred to new tasks.

### 3.2. Single-center loss(SCL)

In this paper, we adopt the idea of single-center loss [25] to further improve the model's generalization capability. Fig. 2 illustrates the of feature distribution w/o and w/ SCL, where circles with different colors indicate different manipulation methods while triangles represent real samples. SCL is designed to make the feature distribution of real faces more compact and, at the same time, move fake features away from the center of real features (red triangle in Fig. 2 (b)). Eqn. (5) shows the single-center loss func-

tion:

$$L_{SCL} = d_{real} + max(d_{real} - d_{fake} + margin, 0) \quad (5)$$

where $d$ is the average distance between the real center and each feature, as shown in Eqn. (6):

$$d = \frac{1}{N} \sum_{i=1}^{N} ||f_i - C||_2 \quad (6)$$

where $C$ represents the real center. We pick the features after the second last fully-connected layer to calculate $L_{SCL}$. By using such loss, the features of real and fake faces become more discriminative and separable, thus leading to a more general face forgery detection performance.

## 4. Experiments

In this section, we conduct cross-domain experiments, including cross-manipulation and cross-dataset settings, to examine the robustness of the model. Then, we perform ablation studies to demonstrate the effectiveness of the LoRA module and the single-center loss.

### 4.1 Implementation details

The proposed framework is implemented by Pytorch. The model is trained using Adam optimizer with $\beta_1$=0.9 and $\beta_2$=0.999. We set the learning rate and weight decay as 1e-4 and 1e-5, respectively. The model is trained on 1 RTX 2080Ti GPUs with batch size 36. The FaceForensics++[35] dataset is used as our training set. We follow the data split strategy in LTW [37] for fair comparison.

### 4.2 Evaluation on cross-manipulation detection

FaceForensics++[35] dataset provides 4000 fake videos generated by four manipulation techniques: Deepfake(DF),

Table 2: Low quality (c40) cross-manipulation detection performance on unseen forgery methods.

| Setting | FF,FS,NT→DF | | DF,FS,NT→FF | | DF,FF,NT→FS | | DF,FF,FS→NT | | Average | |
| Method | AUC | ACC | AUC | ACC | AUC | ACC | AUC | ACC | AUC | ACC |
|---|---|---|---|---|---|---|---|---|---|---|
| ResNet18 [14] | 0.741 | 0.673 | 0.648 | 0.600 | 0.634 | 0.594 | 0.598 | 0.567 | 0.655 | 0.609 |
| Xception [8] | 0.766 | 0.694 | 0.696 | 0.643 | 0.626 | 0.593 | 0.597 | 0.552 | 0.671 | 0.621 |
| EfficientNet [38] | 0.451 | 0.485 | 0.537 | 0.505 | 0.512 | 0.503 | 0.499 | 0.497 | 0.500 | 0.498 |
| All-train EfficientNet [38] | 0.753 | 0.676 | 0.674 | 0.614 | 0.614 | 0.580 | 0.600 | 0.564 | 0.660 | 0.609 |
| Focal-loss EfficientNet [38] | 0.749 | 0.674 | 0.672 | 0.610 | 0.596 | 0.575 | 0.605 | 0.566 | 0.656 | 0.606 |
| Forensics Transfer [9] | N.A. | 0.682 | N.A. | 0.550 | N.A. | 0.530 | N.A. | 0.550 | N.A. | 0.578 |
| Multi-task [33] | N.A. | 0.667 | N.A. | 0.565 | N.A. | 0.517 | N.A. | 0.560 | N.A. | 0.577 |
| MLDG [24] | 0.730 | 0.671 | 0.617 | 0.581 | 0.617 | 0.581 | 0.607 | 0.569 | 0.643 | 0.601 |
| LTW [37] | 0.756 | 0.691 | **0.724** | **0.657** | 0.681 | 0.625 | **0.608** | **0.585** | 0.692 | 0.640 |
| ViT base [12] | 0.739 | 0.643 | 0.650 | 0.595 | 0.592 | 0.560 | 0.552 | 0.538 | 0.633 | 0.584 |
| Ours | **0.818** | **0.735** | 0.686 | 0.638 | **0.710** | **0.653** | 0.582 | 0.542 | **0.699** | **0.642** |

Face2Face(FF), FaceSwap(FS), and NeuralTextures(NT). Each video has three compression levels with different QPs: raw(QP=0), high quality(HQ QP=23), and low quality(LQ QP=40). To examine the generalization capability of the designed model on unseen manipulation techniques and accommodate real-world application scenarios, we conduct cross-manipulation evaluations on both HQ and LQ data, introduced next.

**Detection results on HQ data.** We apply leave-one-out cross-validation and average the results of four trials. Following prior arts, we report AUC and ACC scores in Table 1. Compared with the baseline method (ViT base), the proposed method demonstrated a significant improvement over the baseline method (ViT base), with the average AUC score increasing from 0.623 to 0.792. This improvement can be attributed to the effectiveness of LoRA and SCL. On the other hand, our method is superior to the SOTA method LTW[37] in terms of the average detection performance.

**Detection results on LQ data.** Detecting low-quality manipulated faces is more challenging because severe compression can erase abundant forgery cues. Table 2 presents the detection results on the low-quality data. Compared to the ViT baseline, the average AUC and ACC scores of the proposed method get significant improvements: 6.6% and 5.8%, respectively. Additionally, our model achieves the best average detection performance, demonstrating its outstanding robustness under such a challenging setting.

### 4.3 Evaluation on cross-dataset detection

Evaluating the model on an unseen dataset is another practical scenario where the detection performance of most methods tends to degrade dramatically due to domain shift. In this paper, we train our model on FF++ Deepfake(both c23 and c40) subset and test it on unseen Deepfake datasets, including CelebDF [29], DFD [4], DFDC [11], and Deepfake-TIMIT [23]. The AUC detection scores are reported in Table 3, we can readily observe that the proposed method obtains 11.19% AUC boost compared to the ViT baseline, demonstrating our model's generalization capability from another point of view.

### 4.4 Ablation study

To validate the effectiveness of LoRA and SCL in the task of general face forgery detection, we conduct an ablation evaluation under the challenging LQ cross-manipulation setting. As shown in Table 4, the usage of LoRA modules significantly improves the AUC and ACC scores, and the SCL further boosts the low-quality cross-manipulation detection performance.

## 5. Conclusions

In this paper, we presented a general and robust face forgery detection method based on ViT backbone. Firstly, the backbone is initialized with ImageNet weights, and the loaded parameters are frozen during the training process. Then, we tune the LoRA modules under the joint supervision of cross-entropy and single center losses. By doing this, the number of trainable parameters can be greatly reduced and much computational resource can be saved. Extensive experiments demonstrate that the use of LoRA and SCL can improve the generalization capability of the forgery detection model. The proposed method can serve as a basis for developing ViT-based face forgery detection models.

## 6 Acknowledgement

Table 3: Cross-dataset evaluation results.

| Dataset | CelebDF | DFD (HQ) | DFD (LQ) | DFDC | DFMIT (HQ) | DFMIT (LQ) | Average |
|---|---|---|---|---|---|---|---|
| MesoNet [5] | 58.85 | 62.07 | 52.25 | 54.60 | 33.61 | 45.08 | 51.08 |
| MesoIncep4 [5] | 68.26 | 79.18 | 63.27 | 61.92 | 16.12 | 27.47 | 52.70 |
| ResNet50 [14] | 67.09 | 69.60 | 60.61 | 61.97 | 41.95 | 47.27 | 58.08 |
| Face X-ray [26] | 71.89 | 69.61 | 62.89 | 58.97 | 42.52 | 50.05 | 59.32 |
| DFFD [10] | 69.55 | 71.69 | 60.60 | 59.72 | 32.91 | 39.32 | 55.63 |
| Multi-task [33] | 65.18 | 70.75 | 58.61 | 57.38 | 16.53 | 15.59 | 47.34 |
| EfficientNet [38] | 75.90 | 80.63 | 64.19 | 66.39 | 29.12 | 28.34 | 57.43 |
| $F^3$Net [34] | 72.28 | 72.92 | 58.89 | 63.33 | 38.55 | 45.67 | 58.61 |
| Xception [8] | 67.75 | 72.45 | 59.73 | 63.12 | 33.82 | 40.79 | 56.28 |
| D&L (Effi.) [18] | 67.15 | 73.52 | 67.21 | 60.32 | 49.90 | 51.01 | 61.52 |
| D&L (Xcep.) [18] | 70.65 | 76.23 | 64.53 | 63.31 | 47.20 | **56.08** | 63.00 |
| ViT-base [12] | 71.23 | 61.32 | 59.51 | 66.84 | 56.69 | 47.41 | 60.50 |
| Ours | **83.76** | **83.42** | **68.31** | **71.74** | **70.36** | 52.52 | **71.69** |

Table 4: Ablation study (c40 cross-manipulation).

| ViT | LoRA | SCL | AUC | ACC |
|---|---|---|---|---|
| √ | - | - | 0.633 | 0.584 |
| √ | √ | - | 0.684 | 0.630 |
| √ | √ | √ | **0.699** | **0.642** |

# References

[1] Fakeapp. [EB/OL], 2018. https://www.malavida.com/en/soft.

[2] Github deepfake faceswap. https://github.com/deepfakes/faceswap, 2018.

[3] Deepfakes faceswap. [EB/OL], 2019. https://github.com/deepfakes/faceswap.

[4] Deepfakedetection. https://ai.googleblog.com/2019/09/contributing-data-to-deepfake-detection.html, 2020.

[5] 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)*, pages 1–7. IEEE, 2018.

[6] R. Cai, H. Li, S. Wang, C. Chen, and A. C. Kot. Drl-fas: A novel framework based on deep reinforcement learning for face anti-spoofing. *IEEE Transactions on Information Forensics and Security*, 16:937–951, 2020.

[7] R. Cai, Z. Li, R. Wan, H. Li, Y. Hu, and A. C. Kot. Learning meta pattern for face anti-spoofing. *IEEE Transactions on Information Forensics and Security*, 17:1201–1213, 2022.

[8] F. Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017.

[9] D. Cozzolino, J. Thies, A. Rössler, C. Riess, M. Nießner, and L. Verdoliva. Forensictransfer: Weakly-supervised domain adaptation for forgery detection. *arXiv preprint arXiv:1812.02510*, 2018.

[10] H. Dang, F. Liu, J. Stehouwer, X. Liu, and A. K. Jain. On the detection of digital face manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5781–5790, 2020.

[11] B. Dolhansky, J. Bitton, B. Pflaum, J. Lu, R. Howes, M. Wang, and C. Canton Ferrer. The deepfake detection challenge dataset. *arXiv e-prints*, pages arXiv–2006, 2020.

[12] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

[13] A. Haliassos, R. Mira, S. Petridis, and M. Pantic. Leveraging real talking faces via self-supervision for robust forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14950–14962, 2022.

[14] 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*, pages 770–778, 2016.

[15] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR, 2019.

[16] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.

[17] M. Jia, L. Tang, B.-C. Chen, C. Cardie, S. Belongie, B. Hariharan, and S.-N. Lim. Visual prompt tuning. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*, pages 709–727. Springer, 2022.

[18] C. Kong, B. Chen, H. Li, S. Wang, A. Rocha, and S. Kwong. Detect and locate: Exposing face manipulation by semantic- and noise-level telltales. *IEEE Transactions on Information Forensics and Security*, 17:1741–1756, 2022.

[19] C. Kong, B. Chen, W. Yang, H. Li, P. Chen, and S. Wang. Appearance matters, so does audio: Revealing the hidden face via cross-modality transfer. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(1):423–436, 2021.

[20] C. Kong, S. Wang, and H. Li. Digital and physical face attacks: Reviewing and one step further. *arXiv preprint arXiv:2209.14692*, 2022.

[21] C. Kong, K. Zheng, Y. Liu, S. Wang, A. Rocha, and H. Li. M3fas: An accurate and robust multimodal mobile face anti-spoofing system. *arXiv preprint arXiv:2301.12831*, 2023.

[22] C. Kong, K. Zheng, S. Wang, A. Rocha, and H. Li. Beyond the pixel world: A novel acoustic-based face anti-spoofing system for smartphones. *IEEE Transactions on Information Forensics and Security*, 17:3238–3253, 2022.

[23] P. Korshunov and S. Marcel. Deepfakes: a new threat to face recognition? assessment and detection. *arXiv preprint arXiv:1812.08685*, 2018.

[24] 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*, volume 32, 2018.

[25] J. Li, H. Xie, J. Li, Z. Wang, and Y. Zhang. Frequency-aware discriminative feature learning supervised by single-center loss for face forgery detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6458–6467, 2021.

[26] 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*, pages 5001–5010, 2020.

[27] Y. Li, M.-C. Chang, and S. Lyu. In ictu oculi: Exposing ai created fake videos by detecting eye blinking. In *2018 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–7. IEEE, 2018.

[28] Y. Li and S. Lyu. Exposing deepfake videos by detecting face warping artifacts. *arXiv preprint arXiv:1811.00656*, 2018.

[29] 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*, pages 3207–3216, 2020.

[30] Z. Li, R. Cai, H. Li, K.-Y. Lam, Y. Hu, and A. C. Kot. One-class knowledge distillation for face presentation attack detection. *IEEE Transactions on Information Forensics and Security*, 17:2137–2150, 2022.

[31] 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*, pages 16317–16326, 2021.

[32] I. Masi, A. Killekar, R. M. Mascarenhas, S. P. Gurudatt, and W. AbdAlmageed. Two-branch recurrent network for isolating deepfakes in videos. In *European Conference on Computer Vision*, pages 667–684. Springer, 2020.

[33] 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)*, pages 1–8. IEEE, 2019.

[34] 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*, pages 86–103. Springer, 2020.

[35] 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*, pages 1–11, 2019.

[36] K. Shiohara and T. Yamasaki. Detecting deepfakes with self-blended images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18720–18729, 2022.

[37] 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*, volume 35, pages 2638–2646, 2021.

[38] M. Tan and Q. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.

[39] J. Thies, M. Zollhöfer, and M. Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019.

[40] 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*, pages 2387–2395, 2016.

[41] X. Yang, Y. Li, and S. Lyu. Exposing deep fakes using inconsistent head poses. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8261–8265. IEEE, 2019.

[42] Y. Yu, Y. Wang, W. Yang, S. Lu, Y.-p. Tan, and A. C. Kot. Backdoor attacks against deep image compression via adaptive frequency trigger. *arXiv preprint arXiv:2302.14677*, 2023.

[43] Y. Yu, W. Yang, Y.-P. Tan, and A. C. Kot. Towards robust rain removal against adversarial attacks: A comprehensive benchmark analysis and beyond. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6013–6022, 2022.

[44] Y. Zhou and S.-N. Lim. Joint audio-visual deepfake detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14800–14809, 2021.