



# Exploring Frequency Adversarial Attacks for Face Forgery Detection

Shuai Jia<sup>1</sup> Chao Ma<sup>1\*</sup> Taiping Yao<sup>2</sup> Bangjie Yin<sup>2</sup> Shouhong Ding<sup>2</sup> Xiaokang Yang<sup>1</sup>

<sup>1</sup>MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University

<sup>2</sup>Youtu Lab, Tencent

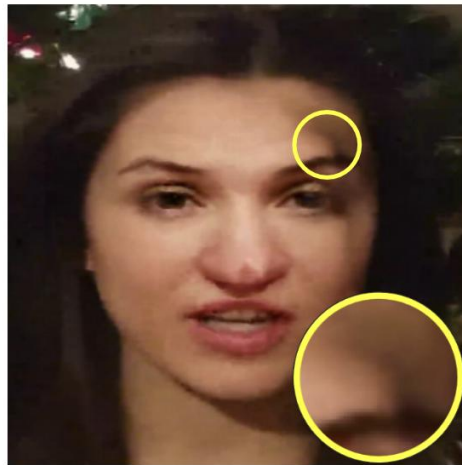
CVPR 2022



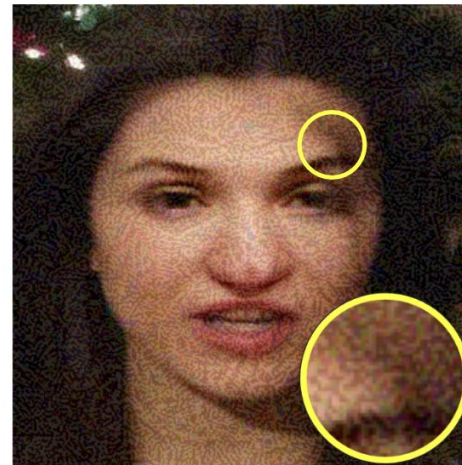


# Motivation

- Various **facial manipulation** techniques have drawn serious public concerns in **morality, security, and privacy**.
- Existing face forgery classifiers are **vulnerable** to adversarial examples with injected **imperceptible perturbations** on the pixels.
- Many face forgery detectors always utilize the **frequency diversity** between real and fake faces as a crucial clue.



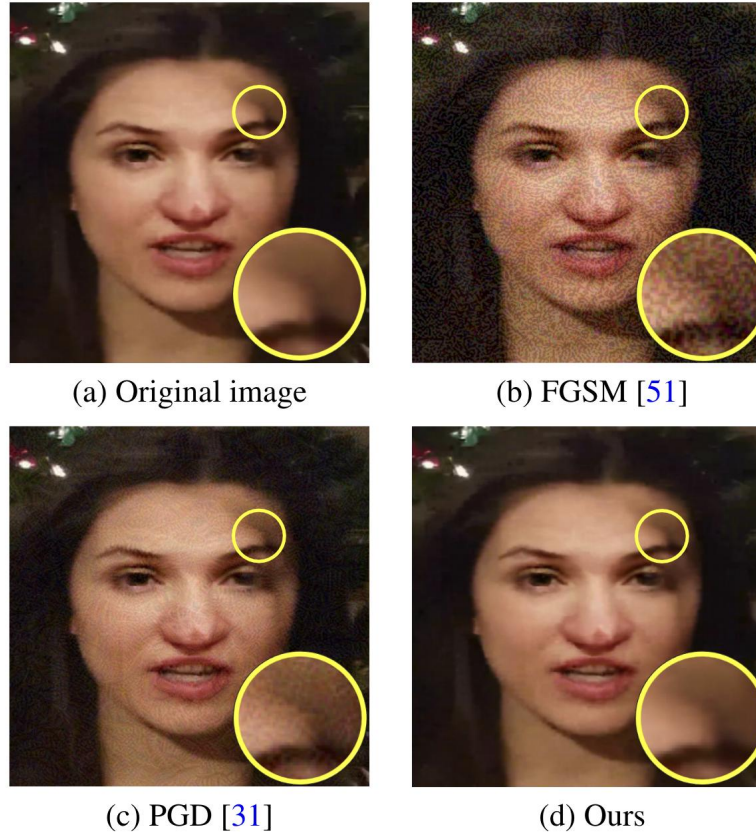
(a) Original image



(b) FGSM [51]



# Motivation



In this paper, instead of injecting adversarial perturbations into the spatial domain, we propose a **frequency adversarial attack method** against face forgery detectors.



# Contributions

- (a) For the task of face forgery detection, we propose a novel adversarial attack method to **generate perturbations in the frequency domain**. Compared with the previous attacks, our method generates more imperceptible perturbations for human observers.
- (b) To further boost the transferability of the attack, we propose a **hybrid adversarial attack** based on the strategy of meta-learning to simultaneously perform attacks on the spatial and frequency domain.
- (c) We perform the proposed method both on the spatial-based face forgery detectors and the state-of-the-art frequency-based detectors. Extensive experiments on benchmarks demonstrate the **effectiveness** of our attack under both white-box and black-box settings.

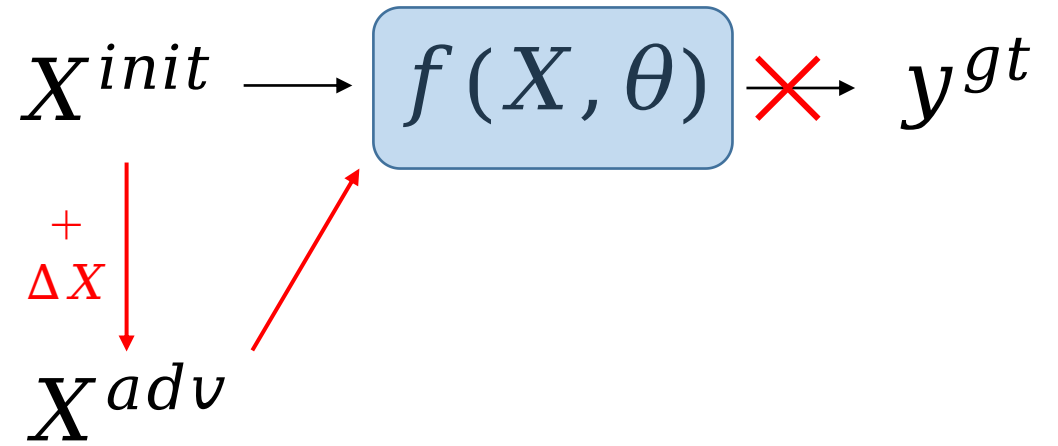


# Methodology

Let  $X^{\text{init}}$  denote the original image,  $f(X, \theta)$  denote the face forgery detector, and  $y^{\text{gt}}$  denote the corresponding ground-truth label. Our aim is to generate the adversarial example  $X^{\text{adv}}$  that makes the face forgery detector predict wrongly, i.e.,  $f(X^{\text{adv}}, \theta) \neq y^{\text{gt}}$ . During adversarial attack, the objective is to maximize the loss function  $\mathcal{L}(X^{\text{adv}}, y^{\text{gt}})$ , where  $\mathcal{L}$  is the binary cross entropy loss in face forgery detection. The concrete optimization is defined as:

$$\arg \max \mathcal{L}(X^{\text{adv}}, y^{\text{gt}}), \quad \text{s.t. } \|X^{\text{adv}} - X^{\text{init}}\|_p < \epsilon, \quad (1)$$

where  $p$  is  $l_p$ -norm to ensure the adversarial image close to the original image.







# Spatial Adversarial Attack

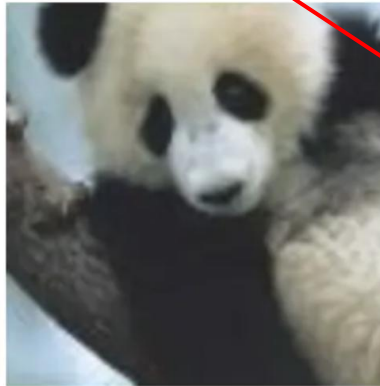
- Fast Gradient Sign Method (FGSM). FGSM is a **single-step** attack method that calculates the perturbations based on the gradient of the adversarial loss. The optimization is defined as:

$$\underline{X^{\text{adv}}} = \underline{X^{\text{init}}} + \underline{\epsilon} \cdot \underline{\text{sign}(\nabla_X \mathcal{L}(X^{\text{adv}}, y^{\text{gt}}))}. \quad (2)$$



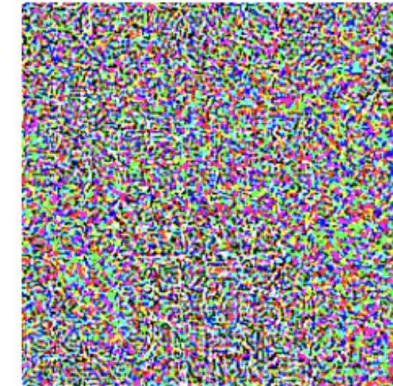
“gibbon”  
99.3% confidence

=



“panda”  
57.7% confidence

+ .007 ×



“nematode”  
8.2% confidence



# Spatial Adversarial Attack

- Fast Gradient Sign Method (FGSM). FGSM is a **single-step** attack method that calculates the perturbations based on the gradient of the adversarial loss. The optimization is defined as:

$$X^{\text{adv}} = X^{\text{init}} + \epsilon \cdot \text{sign}(\nabla_X \mathcal{L}(X^{\text{adv}}, y^{\text{gt}})). \quad (2)$$

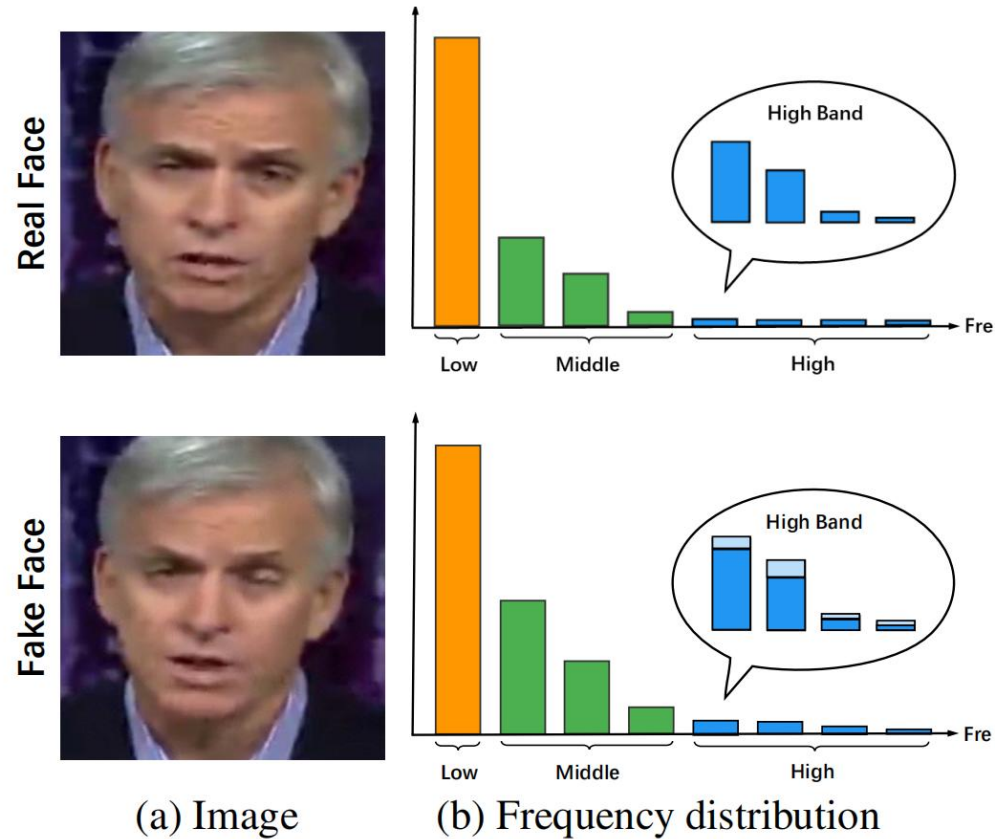
- Projected Gradient Descent (PGD). PGD is a **multi-step** variant of FGSM. Meanwhile, it adopts a random initialization of perturbations at the first step. The update procedure is defined as:

$$\begin{aligned} X_0^{\text{adv}} &= X^{\text{init}}, \\ X_{n+1}^{\text{adv}} &= \text{Clip} \left\{ X_n^{\text{adv}} + \alpha \cdot \text{sign}(\nabla_X \mathcal{L}(X_n^{\text{adv}}, y^{\text{gt}})) \right\}. \end{aligned} \quad (3)$$



# Frequency Adversarial Attack

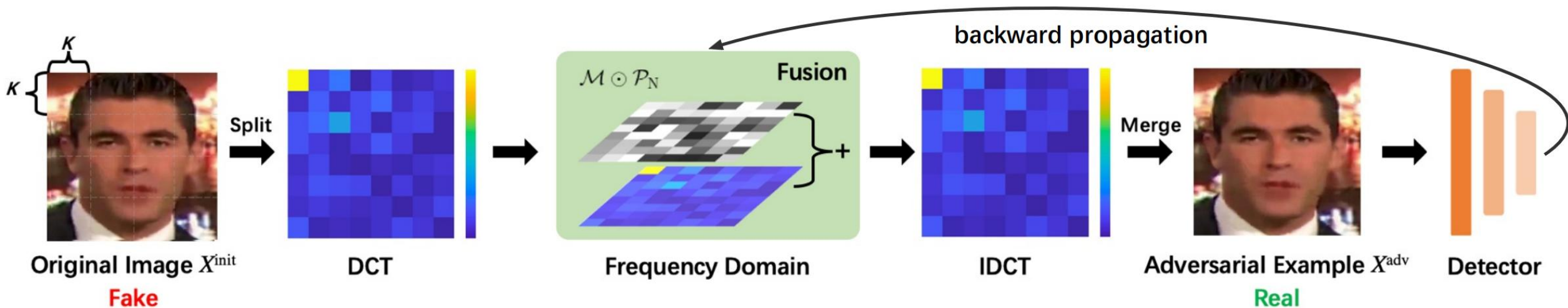
Previous studies<sup>[1][2]</sup> have proven the difference between the real face and the fake face in the frequency domain.



- [1] Shen Chen, Taiping Yao, Yang Chen, Shouhong Ding, JilinLi, and Rongrong Ji. Local relation learning for face forgery detection. In AAIL, 2021.
- [2] Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and JingShao. Thinking in frequency: Face forgery detection by min-ing frequency-aware clues. In ECCV, 2020.



# Frequency Adversarial Attack

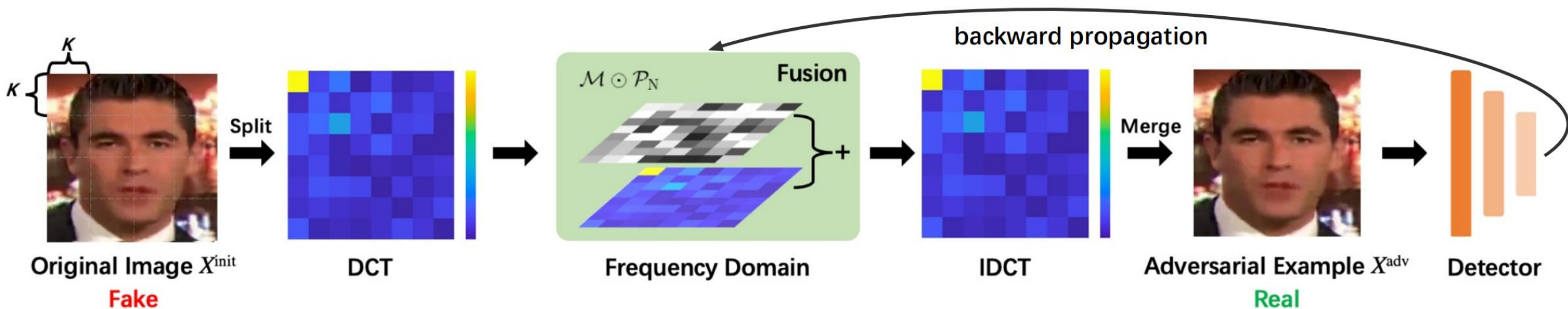


We summarize the optimization procedure as follows:

$$\begin{aligned} \arg \max \quad & \mathcal{L}(\mathcal{D}'(\mathcal{F}(\mathcal{D}(X^{\text{adv}}))), \theta, y^{\text{gt}}), \\ \text{s.t.} \quad & \|\mathcal{D}(X^{\text{adv}}) - \mathcal{D}(X^{\text{init}})\|_{\text{p}} < \epsilon, \end{aligned} \quad (4)$$

where  $\mathcal{D}(\cdot)$  denotes discrete cosine transform (DCT),  $\mathcal{D}'(\cdot)$  denotes inverse discrete cosine transform (IDCT),  $\mathcal{F}$  represents the fusion module to modify the energy in the frequency domain.

# Frequency Adversarial Attack



The complete fusion module is defined as:

$$\mathcal{F}(X_n^{\text{adv}}) = \mathcal{D}(X_n^{\text{adv}}) + \mathcal{M} \odot \mathcal{P}_{n+1}, \quad (6)$$

where  $\odot$  is Hadamard product. During the optimization,  $\mathcal{P}_{n+1}$  is updated as follows:

$$\mathcal{P}_{n+1} = \mathcal{P}_n + \lambda \cdot \text{sign}(\nabla_{\mathcal{P}} \mathcal{L}(\mathcal{D}'(\mathcal{F}(\mathcal{D}(X_n^{\text{adv}}))), \theta, y^{\text{gt}})), \quad (7)$$

where  $\lambda$  is the step size in each iteration.



# Hybrid Adversarial Attack

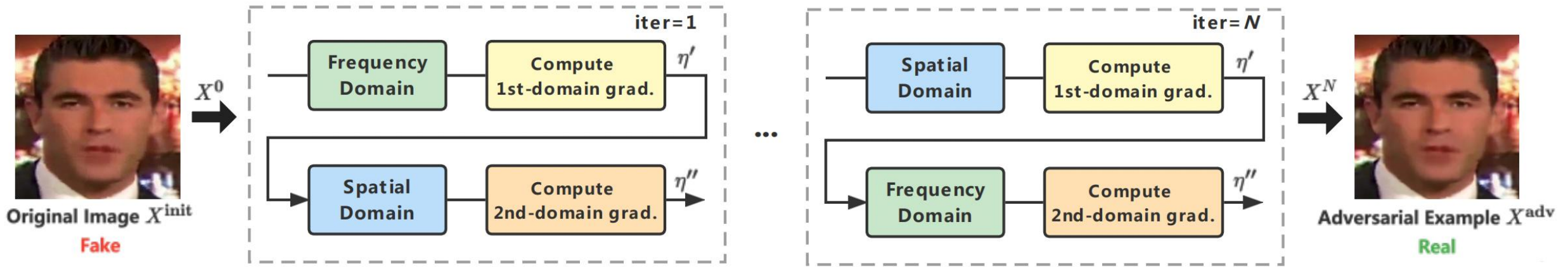


Figure 4. The procedure of hybrid adversarial attack. To combine the adversarial attack in different domains, we calculate gradients from both domains in order and update the perturbations. Then, we switch the order of domains in the next step. After iterations, the adversarial example gathers the gradients from both domains, leading to a stronger adversarial attack on both white-box and black-box settings.



# Experiments

Datasets: DFDC and FaceForensics++

Metric: Attack Success Rate

Forgery Models: Spatial-based: EfficientNet b4, ResNet 50, and XceptionNet;

Frequency-based: F3-Net and LRL.

Table 1. The accuracy of spatial-based and frequency-based face forgery detectors on the DFDC [9] and FaceForensics++ [38] datasets.

Dataset	EfficientNet_b4 [43]	ResNet_50 [20]	XceptionNet [6]	F <sup>3</sup> -Net [5]	LRL [36]
DFDC [9]	91.1%	78.7%	88.0%	69.8%	90.4%
FaceForensics++ [37]	94.3%	89.1%	92.7%	88.8%	98.2%





# Attack on Spatial-based Models

Table 2. The attack success rate of fake faces on spatial-based models on the DFDC [9] dataset.

Model	Attack	Eff_b4 [43]	Res50 [20]	Xcep [6]
Eff_b4 [43]	FGSM	33.2%	7.1%	2.3%
	PGD	77.7%	8.7%	1.8%
	Ours	<b>97.1%</b>	<b>20.1%</b>	<b>2.7%</b>
Res50 [20]	FGSM	0.0%	36.7%	0.9%
	PGD	0.0%	85.4%	0.0%
	Ours	<b>23.2%</b>	<b>87.8%</b>	<b>24.1%</b>
Xcep [6]	FGSM	0.0%	8.4%	45.6%
	PGD	0.0%	10.1%	72.3%
	Ours	<b>1.2%</b>	<b>14.3%</b>	<b>77.5%</b>

Table 3. The attack success rate of fake faces on spatial-based models on the FaceForensics++ [37] dataset.

Model	Attack	Eff_b4 [43]	Res50 [20]	Xcep [6]
Eff_b4 [43]	FGSM	38.7%	4.8%	0.9%
	PGD	71.6%	1.3%	0.3%
	Ours	<b>83.2%</b>	<b>22.7%</b>	<b>1.4%</b>
Res50 [20]	FGSM	3.2%	32.0%	2.1%
	PGD	3.9%	60.2%	2.3%
	Ours	<b>41.4%</b>	<b>65.4%</b>	<b>49.6%</b>
Xcep [6]	FGSM	1.1%	4.1%	18.9%
	PGD	1.1%	7.7%	61.6%
	Ours	<b>1.5%</b>	<b>8.5%</b>	<b>70.5%</b>





# Attack on Frequency-based Models

Table 6. The attack success rate of fake faces on frequency-based models on the DFDC [9] dataset.

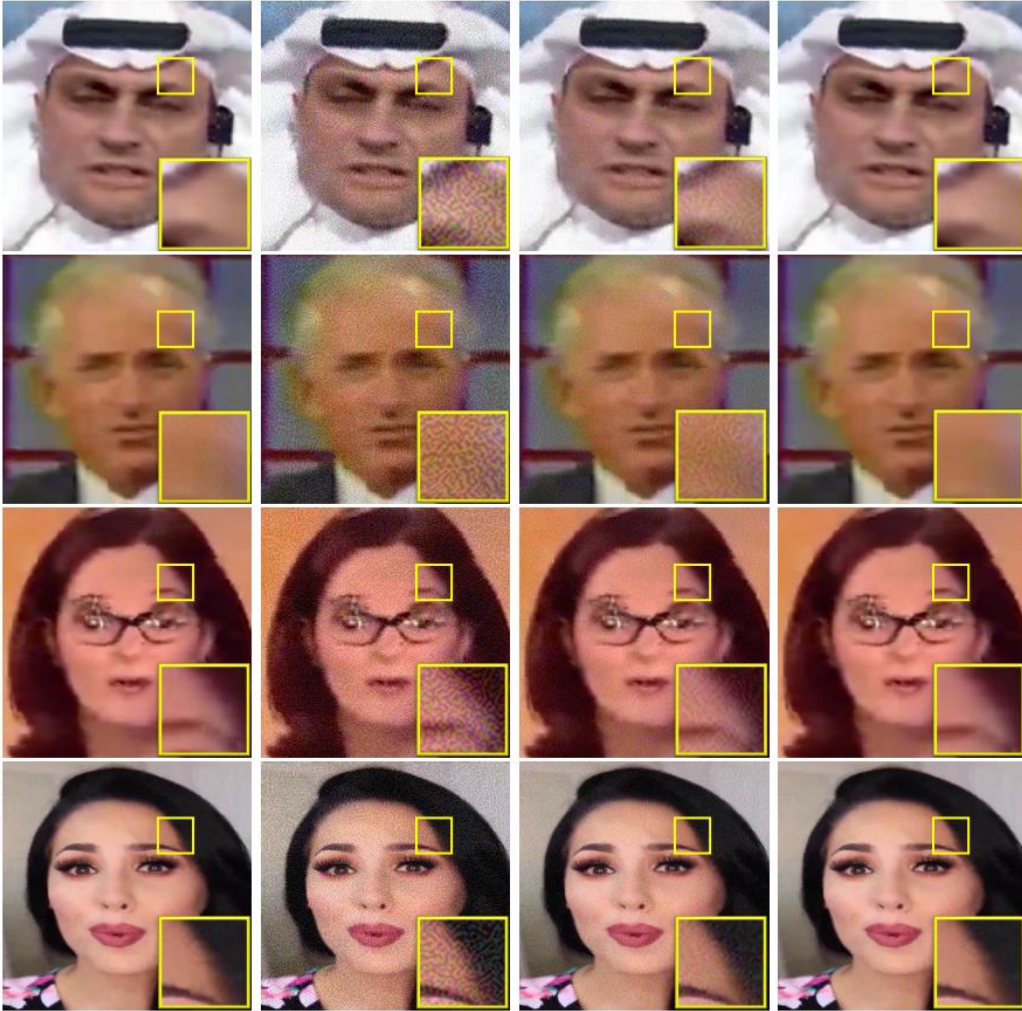
Model	Attack	F <sup>3</sup> -Net [5]	LRL [36]
F <sup>3</sup> -Net [5]	FGSM	43.5%	9.6%
	PGD	97.6%	4.0%
	Ours	<b>98.7%</b>	<b>10.3%</b>
LRL [36]	FGSM	2.3%	71.3%
	PGD	3.0%	<b>100.0%</b>
	Ours	<b>5.5%</b>	<b>100.0%</b>
Eff_b4 [43]	Ours	7.4%	8.5%
Res50 [20]	Ours	<b>12.8%</b>	<b>43.6%</b>
Xcep [6]	Ours	7.6%	9.1%

Table 7. The attack success rate of fake faces on frequency-based models on the FaceForensics++ [37] dataset.

Model	Attack	F <sup>3</sup> -Net [5]	LRL [36]
F <sup>3</sup> -Net [5]	FGSM	24.8%	7.7%
	PGD	80.9%	28.7%
	Ours	<b>82.5%</b>	<b>36.2%</b>
LRL [36]	FGSM	0.2%	68.6%
	PGD	0.0%	98.7%
	Ours	<b>0.5%</b>	<b>99.3%</b>
Eff_b4 [43]	Ours	0.5%	11.8%
Res50 [20]	Ours	<b>7.1%</b>	<b>57.5%</b>
Xcep [6]	Ours	1.1%	19.5%



# Image Quality Assessment



(a) Original (b) FGSM [15] (c) PGD [31] (d) Ours

Table 10. Quantitative evaluation of adversarial examples generated by FGSM [15], PGD [31] and our method on the FaceForensics++ [37] dataset.

Attack method	MSE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
FGSM	0.0279	23.3	0.0881
FGD	0.0238	30.4	0.1343
Ours	<b>0.0027</b>	<b>42.7</b>	<b>0.1763</b>





# Q&A

## Thank you!

