



# Face Forgery Detection by 3D Decomposition

Xiangyu Zhu<sup>1,2</sup>, Hao Wang<sup>1</sup>, Hongyan Fei<sup>3</sup>, Zhen Lei<sup>1,2</sup>, Stan Z. Li<sup>4</sup>

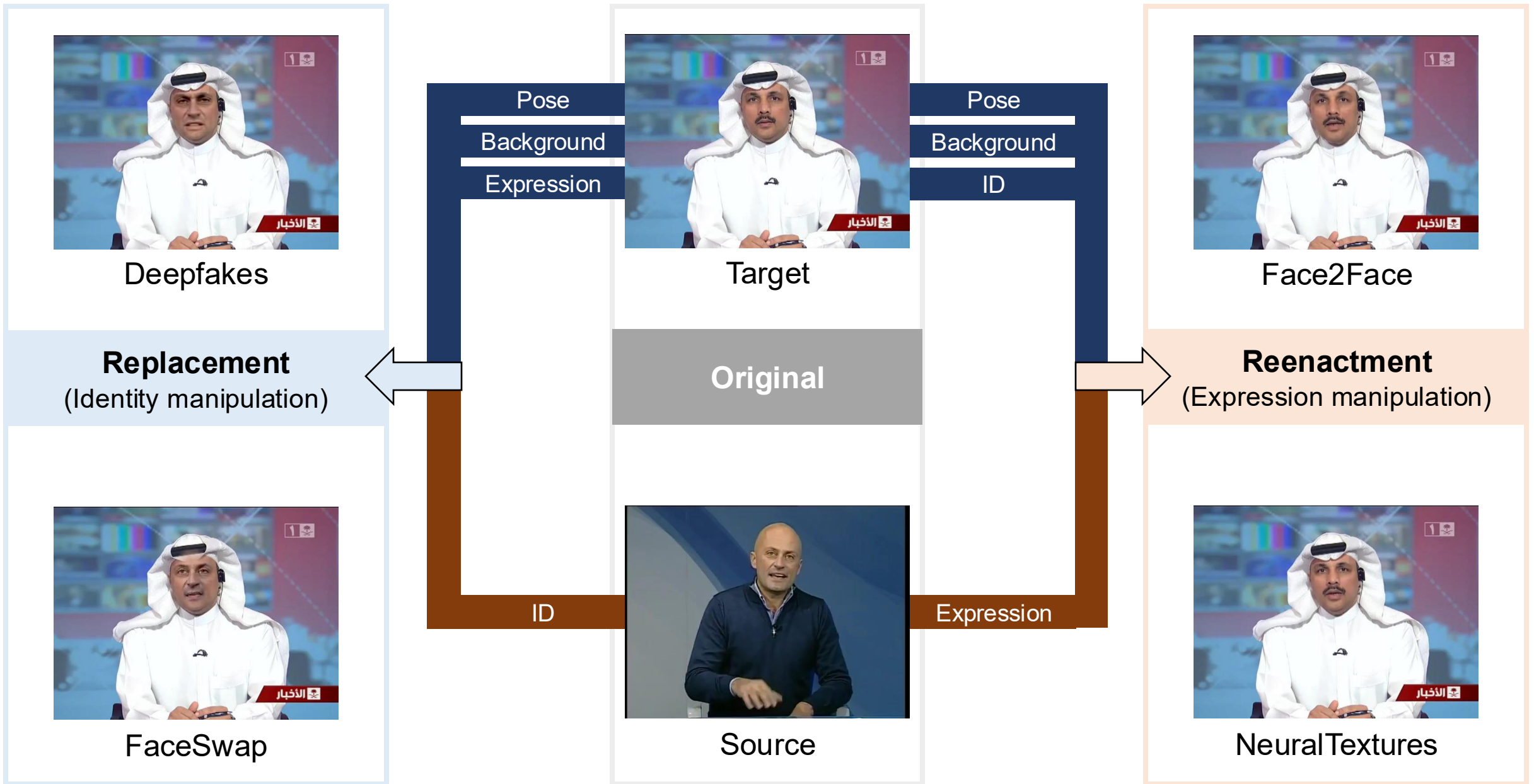
<sup>1</sup>CBSR & NLPR, Institute of Automation, Chinese Academy of Sciences,

<sup>2</sup>School of Artificial Intelligence, University of Chinese Academy of Sciences,

<sup>3</sup>School of Automation and Electrical Engineering, University of Science and Technology Beijing,

<sup>4</sup>School of Engineering, Westlake University

# Face Forgery

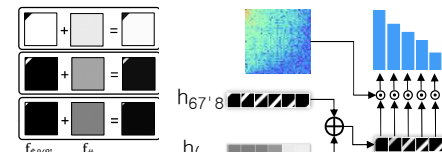


# Goal: explore a novel decomposition

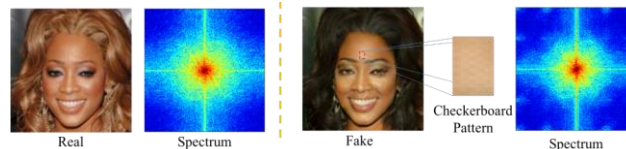
Hand-crafted **frequency** filters  
(Stuchi *et al.* 2017)



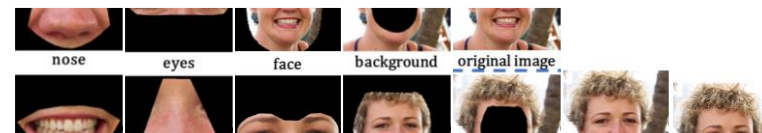
Learned **frequency** filters  
(Qian *et al.* 2020)



Unique replications of **spectra**  
(Zhang *et al.* 2019)

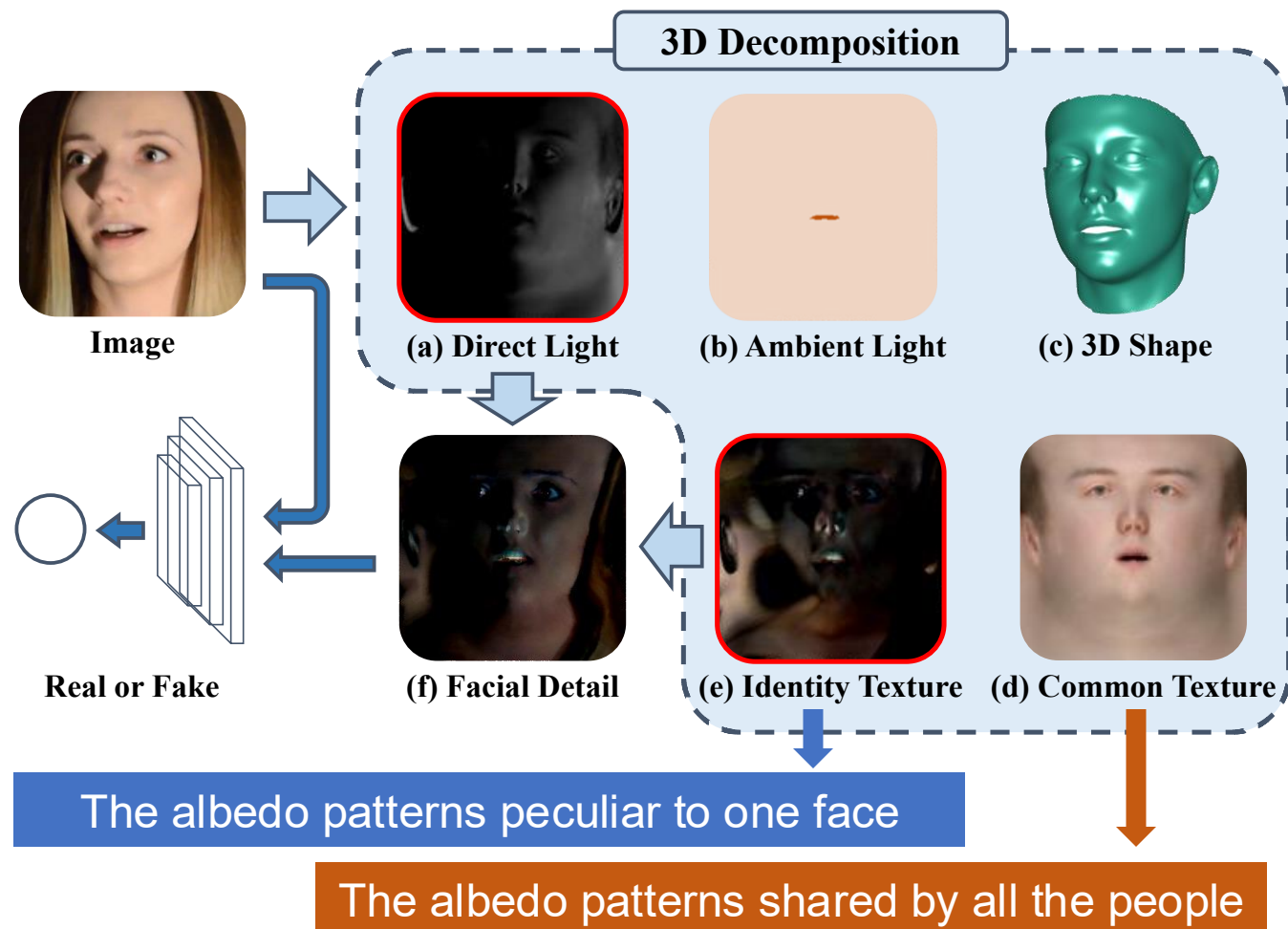


**Spatial-** and **frequency-**domain features  
(Chen *et al.* 2020)



- **Difficult to decide which range of signals contains artifacts**
  - Different capture devices, environments, and compression algorithms
  - Large frequency distribution bias across datasets
- **Suffering from the generalization problem**

# Idea: 3D Decomposition



# 3D Decomposition Implementation

- In computer graphics, a **face image** is generated by:

$$\mathbf{I}_{syn} = Z\text{-Buffer}(\mathbf{S}, \mathbf{C})$$

- Under the *Lambertian assumption*, the RGB of  $i^{\text{th}}$  vertex is:

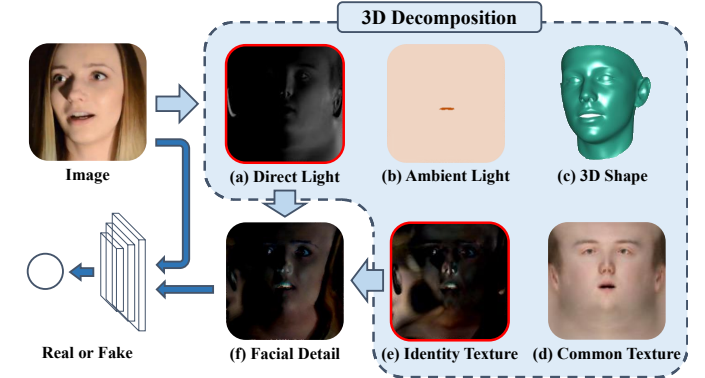
$$\mathbf{C}_i = \mathbf{Amb} * \mathbf{T}_i + \langle \mathbf{n}_i, \mathbf{l} \rangle \cdot \mathbf{Dir} * \mathbf{T}_i$$

- Modeling the common texture by the **BFM PCA texture** model:

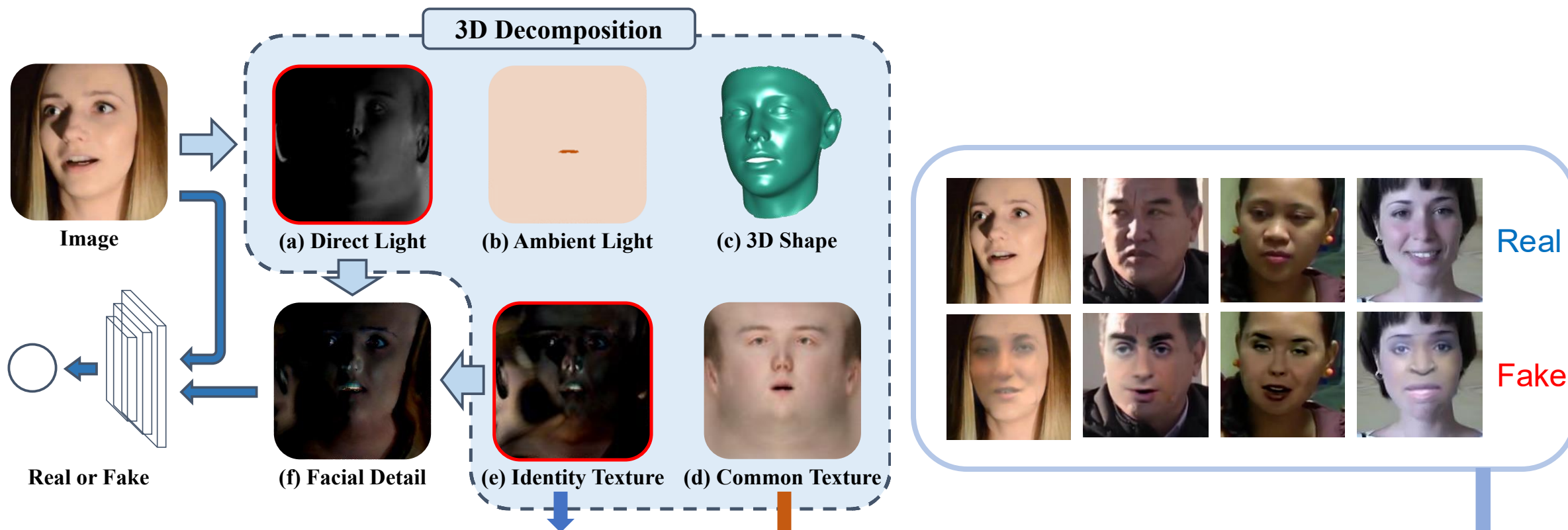
$$\mathbf{T} = \bar{\mathbf{T}} + \mathbf{B}\boldsymbol{\beta} + \mathbf{T}_{id}$$

- Optimizing the following loss to achieve the 3D face decomposition:

$$\arg \min_{\mathbf{S}, \mathbf{Amb}, \mathbf{Dir}, \boldsymbol{\beta}, \mathbf{T}_{id}} \|\mathbf{I} - \mathbf{I}_{syn}(\mathbf{S}, \mathbf{Amb}, \mathbf{Dir}, \boldsymbol{\beta}, \mathbf{T}_{id})\|$$



# Assumption



The albedo patterns peculiar to one face

The albedo patterns shared by all the people

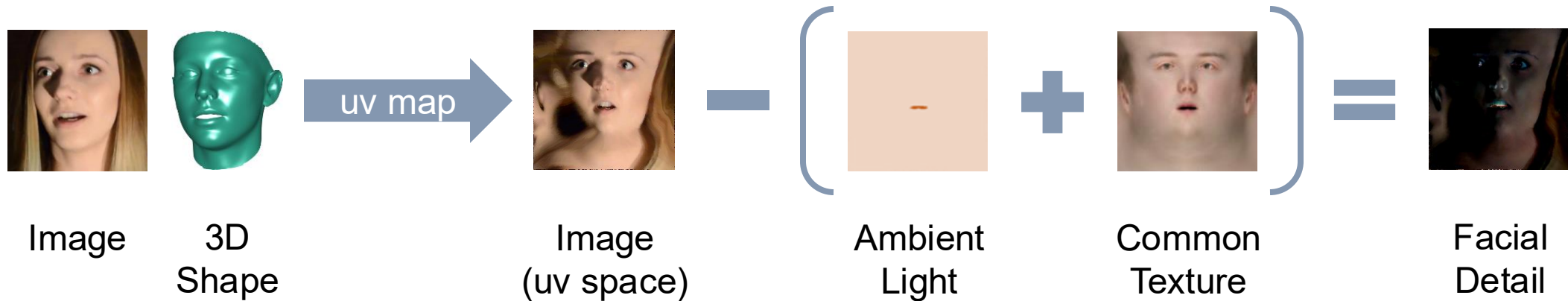
- **Direct light**: decisive forgery clue with the observation on large artifacts under intense direct light
- **Identity texture**: hard to be simulated due to the rich variations across faces
- **3D shape**: its normalizing makes CNN concentrate on specific face regions and simplifies the detection
- **Ambient light and common texture**: easy to be faked and have little forgery clues

- Fitting texture and illumination by **PCA Texture** model and **Spherical Harmonic Reflectance**

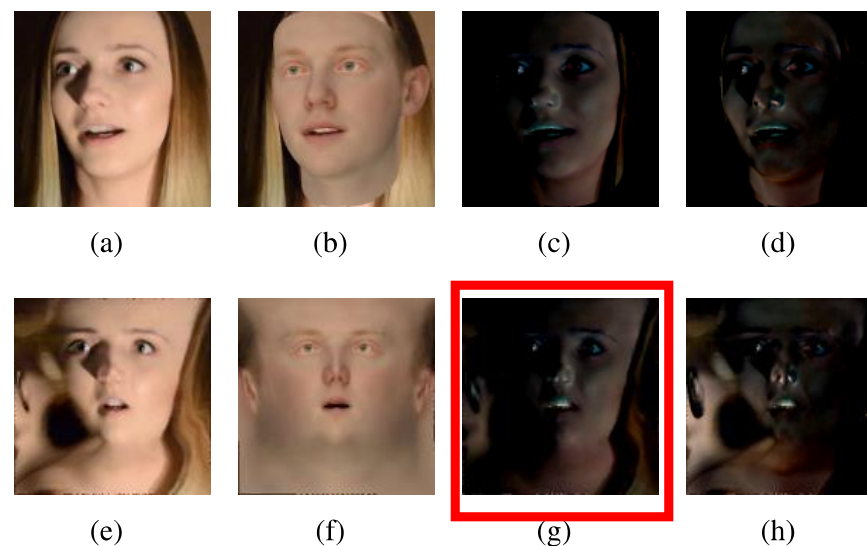
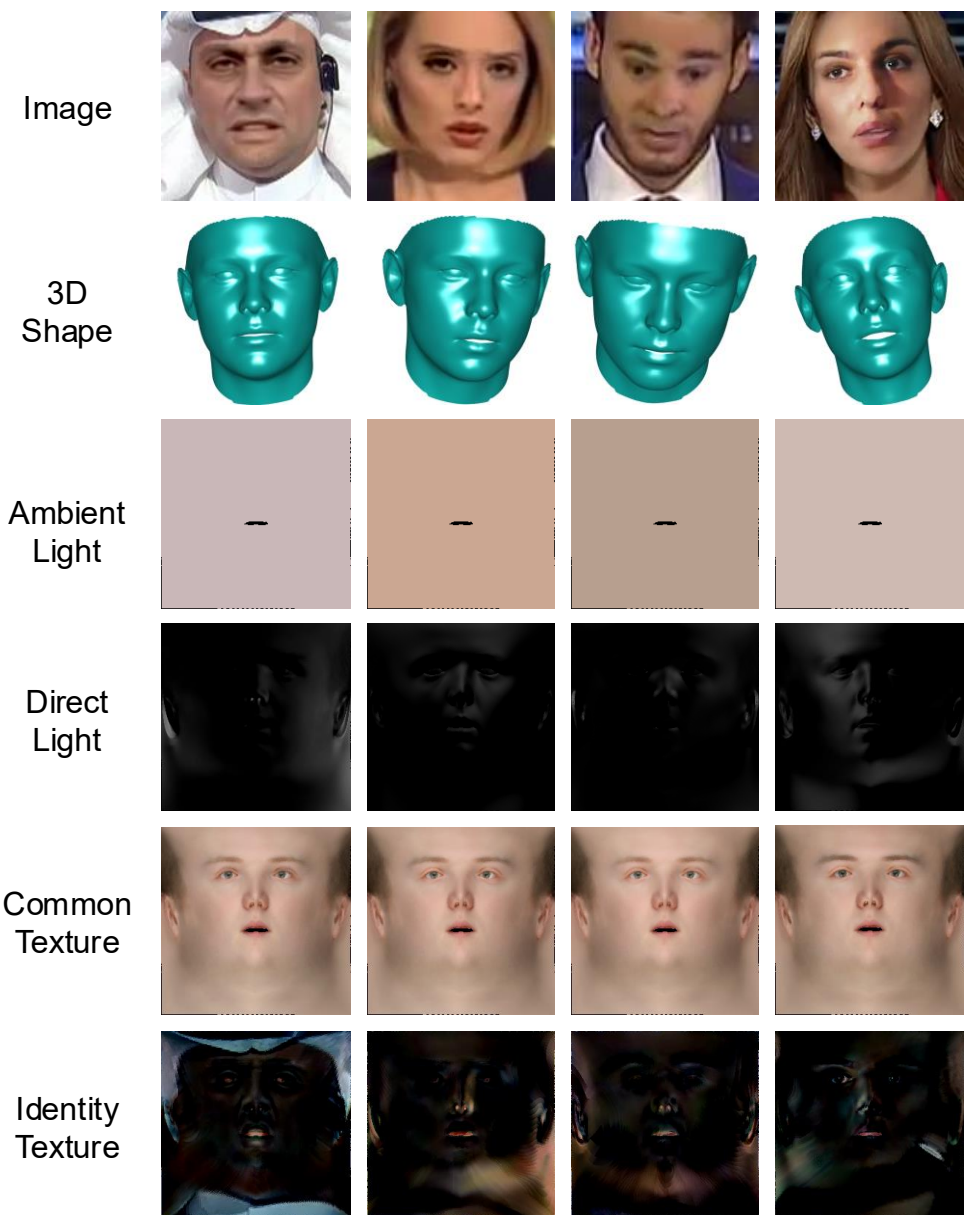
$$\mathbf{I}(\mathbf{S}) = (\mathbf{H}\boldsymbol{\gamma}) \cdot * (\bar{\mathbf{T}} + \mathbf{B}\boldsymbol{\beta})$$

- **Facial detail:** Subtracting face image by the common texture under ambient light (in uv space)

$$\mathbf{FD} = UV(\mathbf{I} - (\mathbf{h}_1\boldsymbol{\gamma}_1) \cdot * (\bar{\mathbf{T}} + \mathbf{B}\boldsymbol{\beta}), \mathbf{S})$$

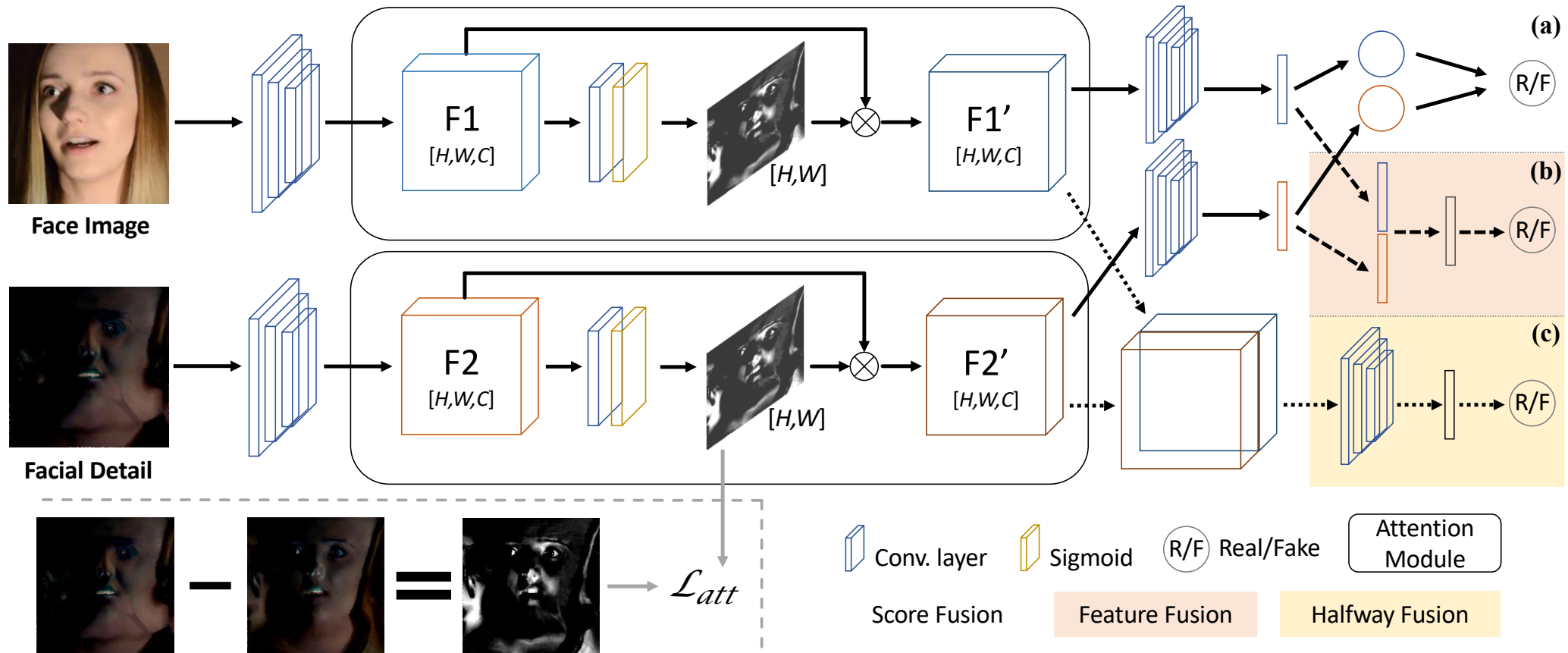


# Component Selection



input	3D	amb	dir	ctex	itex	AUC
In-a	✓	✓	✓	✓	✓	99.13
In-b	✓	✓		✓		50.00
In-c	✓		✓		✓	99.29
In-d	✓				✓	99.14
In-e		✓	✓	✓	✓	98.93
In-f		✓		✓		50.00
<b>In-g</b>			✓		✓	<b>99.56</b>
In-h					✓	99.27

# Forgery Detection with Facial Detail Net



- **Two-stream structure:** regarding **face image** and **facial detail** as two modalities
- **Half-way fusion:** concatenating the **intermediate** 2D feature maps of each stream
- **Detail-guided Attention:** supervised by the **facial detail difference**

# Analysis of the Two-stream Network



Structure	FFpp			DFD			DFDC		
	AP	AUC	EER	AP	AUC	EER	AP	AUC	EER
Img	99.44	99.31	5.39	88.07	65.57	38.38	85.60	62.17	39.99
Detail	99.40	99.12	5.51	87.24	64.29	40.87	85.02	61.80	40.37
Img + Detail (HF)	<b>99.42</b>	<b>98.73</b>	<b>5.63</b>	<b>89.61</b>	<b>78.65</b>	<b>26.03</b>	<b>87.31</b>	<b>66.09</b>	<b>35.46</b>

The halfway fusion makes the fused features correspond to similar receptive fields



Further promotes the results on DFD and DFDC by the local fusion manner

## ➤ Cross-data

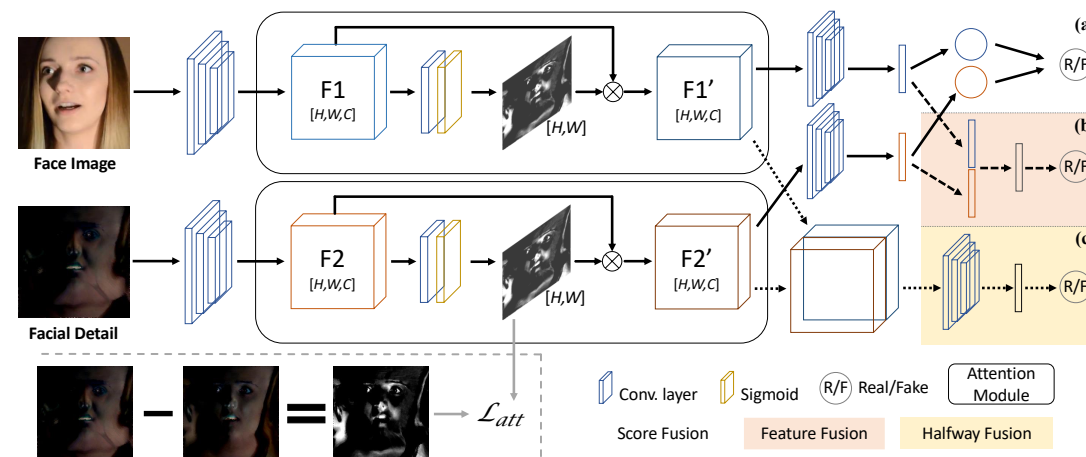
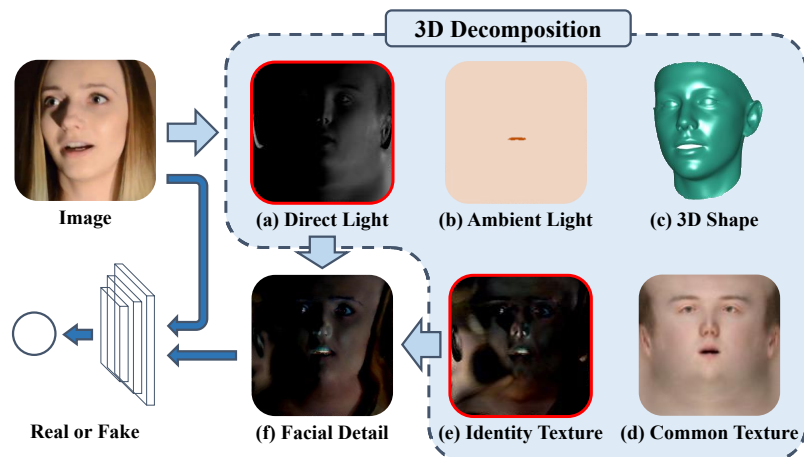
Model	Training dataset	DFD			DFDC		
		AP	AUC	EER	AP	AUC	EER
Xception [50]	FFpp	88.07	65.57	38.38	85.60	62.17	39.99
EfficientNetB4 Ensemble [9]	FFpp	89.35	72.82	34.86	85.71	63.03	38.86
FD <sup>2</sup> Net	FFpp	<b>89.84</b>	<b>79.08</b>	<b>25.18</b>	<b>87.93</b>	<b>67.70</b>	<b>34.91</b>

## ➤ Different Manipulation Methods

Model	Training data	Acc	
		F2F	FS
MesoInception4 [1]	F2F	84.56	56.71
VA-LogReg [39]		83.62	59.45
LAE [19]		90.34	62.51
Multi-task [40]		91.27	55.04
Face X-ray [36]		97.73	85.69
Xception + HP Filter		97.98	57.46
FD <sup>2</sup> Net		<b>98.22</b>	<b>86.54</b>

Effectively extract more discriminative and general features, even from a different distribution of the training dataset

# Conclusion



- A novel face forgery detection method by the **3D decomposition** of the face image
  - 3D shape, common texture, identity texture, ambient light, and direct
  - Critical forgery clues in **direct light** and **identity texture**
  - Propose the **facial detail** to highlight the subtle forgery patterns
- **FD<sup>2</sup>Net**, a two-stream network with facial-detail-supervised attention module
- Improved **effectiveness** and **robustness** of the proposed FD<sup>2</sup>Net
- A **novel direction** to explore the forgery clues by analyzing the physical generation of an image



Thank you!