

R³ Adversarial Network for Cross Model Face Recognition

Ken Chen, Yichao Wu, Haoyu Qin, Ding Liang Xuebo Liu, Junjie Yan

Sensetime Group Limited

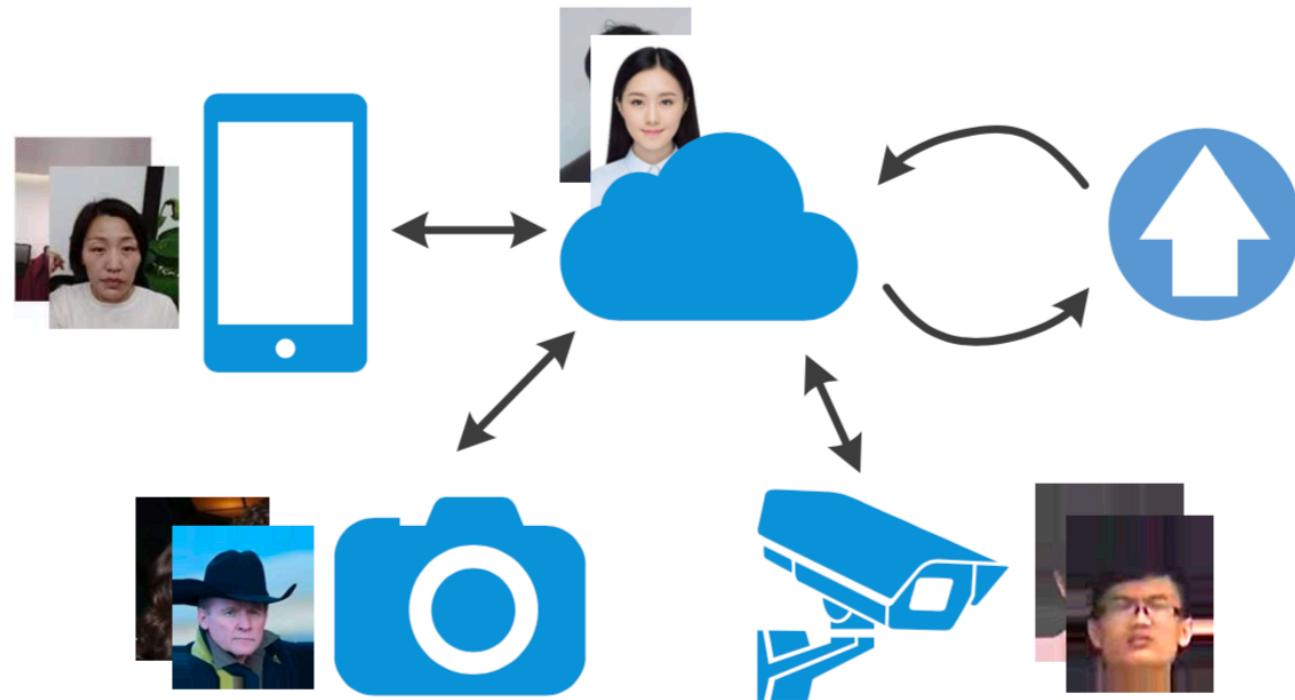
CVPR 2019

Sungman, Cho.

Introduction

Cross Model Face Recognition ?

- Pursuing interaction between information collected from various terminals is a new trend.



Cross Model Face Recognition ?

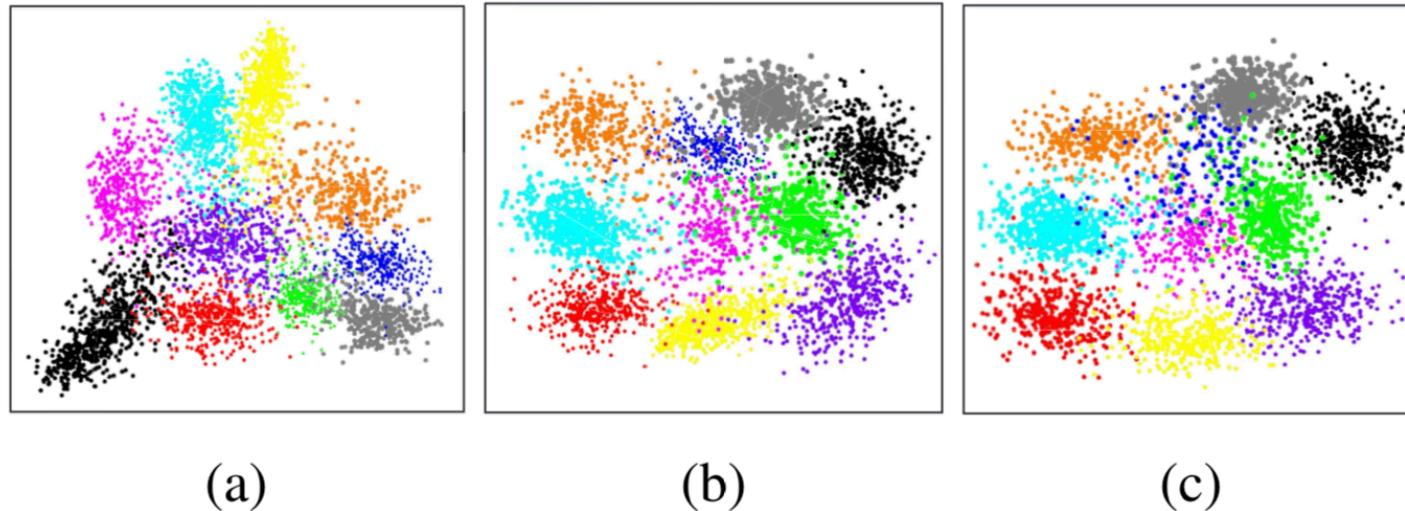


Figure 2: The feature distributions of two typical face recognition models.(a), (b) and (c) are feature distribution of source model, transformation model and target model, respectively.

Cross Model Face Recognition ?

- The core of this problem is to make **features extracted from different models comparable**.
- To solve this problem, from the perspective of Bayesian modeling, we propose R³AN.
- R³AN : Reconstruction, Representation and Regression.
- Introduce **adversarial learning into the reconstruction** path for better performance.

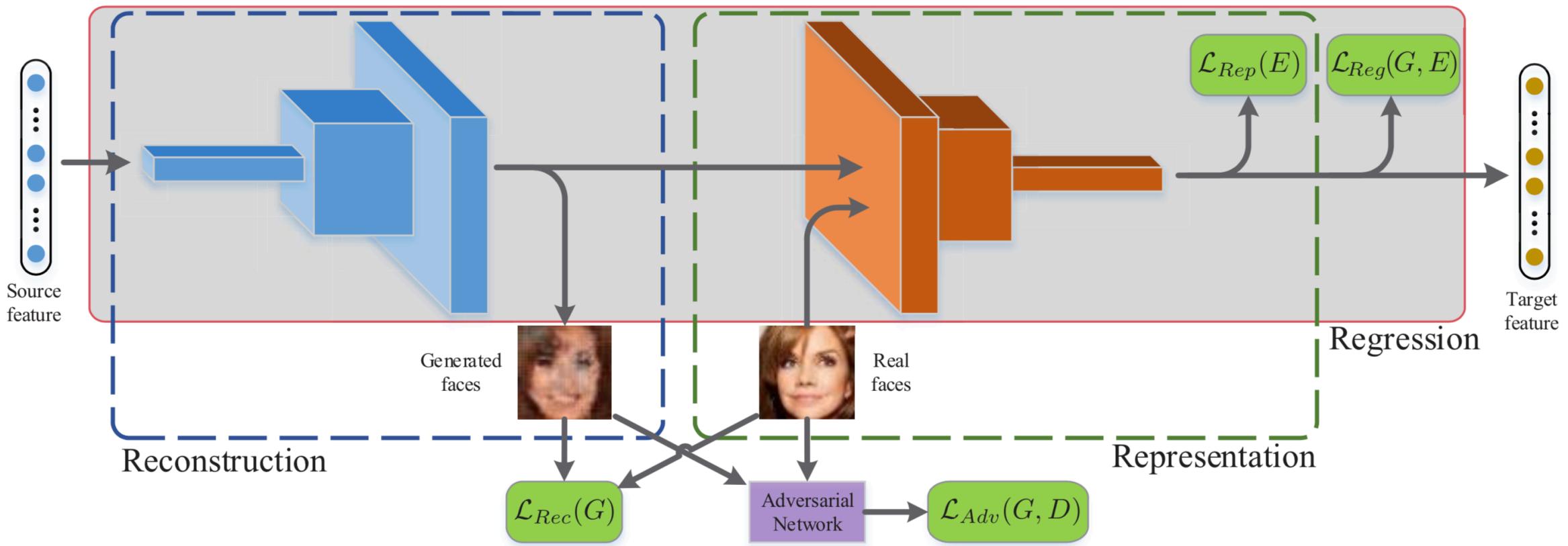
Why R³AN ?

- Single model **lacks the ability** to achieve satisfactory performance when **various domains, applications** and response time requirement are taken into account.
- It violates **privacy policy in common sense**. Uploading and storing user's face images are generally forbidden in industrial community.

Contribution

- Raise the CMFR problem for the first time, which possesses considerable economic and social significance.
- **R³AN** is super fast and valid when solving this problem
- Adversarial learning greatly improves the performance of **R³AN** and recovers higher quality face image, which may **give valuable hints for improving the original face recognition models.**

Architecture



Basic Model

X, Y : source and target feature

Find a mapping function to **maximize the conditional probability $P(Y|X)$.**

map X to Y in one-dimensional space directly by a naïve model like MLP.

→ yields unsatisfactory result.

Bayesian based Model

latent variables $h \in H = \{h_1, \dots, h_K\}$

$$P(\mathbf{Y}|\mathbf{X}) = \sum_{h \in H} P(h|\mathbf{X})P(\mathbf{Y}|\mathbf{X}, h).$$

X and Y is conditionally independent when h is given.

$$P(\mathbf{Y}|\mathbf{X}) = \sum_{h \in H} P(h|\mathbf{X})P(\mathbf{Y}|h).$$

h should be a latent variable independent of models.

$$P(\mathbf{Y}|\mathbf{X}) = P(\mathbf{I}|\mathbf{X})P(\mathbf{Y}|\mathbf{I})$$

Bayesian based Model

latent variables $h \in H = \{h_1, \dots, h_K\}$

$$P(\mathbf{Y}|\mathbf{X}) = \sum_{h \in H} P(h|\mathbf{X})P(\mathbf{Y}|\mathbf{X}, h).$$

X and Y is conditionally independent when h is given.

$$P(\mathbf{Y}|\mathbf{X}) = \sum_{h \in H} P(h|\mathbf{X})P(\mathbf{Y}|h).$$

h should be a latent variable independent of models.

$$P(\mathbf{Y}|\mathbf{X}) = P(\mathbf{I}|\mathbf{X})P(\mathbf{Y}|\mathbf{I})$$

R³AN

- **Reconstruction** : $P(I|X)$ → recover the original image.
- **Representation** : $P(Y|I)$ → extract feature from latent face images
- **Regression** : views reconstruction and representation as a unified problem.
(jointly optimize the whole model)

Reconstruction Path

Restore original face images from extracted features, it can be regarded as generator (G)

Naïve Reconstruction

$$\mathcal{L}_{Rec}(G) = \mathbb{E}_{\mathbf{X}, \mathbf{I}} [||\mathbf{I} - G(\mathbf{X})||_2],$$

Adversarial Reconstruction

$$\mathcal{L}_{Adv}(G, D) = - \mathbb{E}_{\mathbf{X}} [\log (1 - D(G(\mathbf{X})))].$$

Representation Path

Representation path acts as an feature extractor, denoted as E .

It takes the original face images as input and learns representation of the target model.

Representation path can be also considered as a knowledge distiller, which can transfer the knowledge of the target model to the feature extractor module.

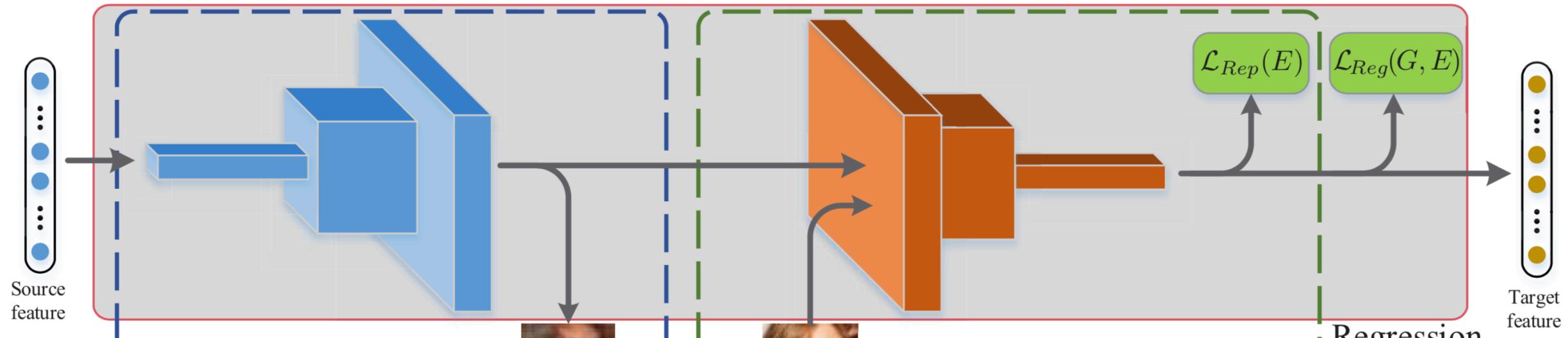
$$\mathcal{L}_{Rep}(E) = \mathbb{E}_{\mathbf{I}, \mathbf{Y}} [||\mathbf{Y} - E(\mathbf{I})||_2].$$

Regression Path

Combines reconstruction and representation together into a unified framework, and is used to jointly optimize the above two path.

$$\mathcal{L}_{Reg}(G, E) = \mathbb{E}_{\mathbf{X}, \mathbf{Y}} [||\mathbf{Y} - E(G(\mathbf{X}))||_2].$$

Architecture



	Input	-	$256 \times 1 \times 1$
Generator	fConv1	ConvTranspose2d	$512 \times 4 \times 4$
	fConv2	ConvTranspose2d	$128 \times 7 \times 7$
	fConv3	ConvTranspose2d	$32 \times 14 \times 14$
	fConv4	ConvTranspose2d	$8 \times 28 \times 28$
	fConv5	ConvTranspose2d	$3 \times 56 \times 56$

	Conv1	Conv2d	$4 \times 28 \times 28$
Discriminator	Conv2	Conv2d	$8 \times 14 \times 14$
	Conv3	Conv2d	$16 \times 7 \times 7$
	Conv4	Conv2d	$32 \times 4 \times 4$
	Conv5	Conv2d	$64 \times 2 \times 2$
	Conv6	Conv2d	$128 \times 1 \times 1$
	FC	Fully Connection	1

Extractor	Conv1	bottleneck	$32 \times 28 \times 28$
	Conv2	bottleneck	$64 \times 14 \times 14$
	Conv3	bottleneck	$96 \times 14 \times 14$
	Conv4	bottleneck	$160 \times 7 \times 7$
	Conv5	bottleneck	$320 \times 7 \times 7$
	Conv6	Conv2d 1×1	$1280 \times 7 \times 7$
	Pooling	AvgPool2d	$1280 \times 1 \times 1$
	FC	Fully Connection	$256 \times 1 \times 1$

Optimization

Final optimization goal :

$$(G^*, E^*) = \arg \min_{G, E} \max_D [\lambda_{Rec} \mathcal{L}_{Rec}(G) + \lambda_{Adv} \mathcal{L}_{Adv}(G, D) \\ + \lambda_{Reg} \mathcal{L}_{Reg}(G, E)] + \lambda_{Rep} \mathcal{L}_{Rep}(E).$$

Optimization

```
1: Random initialize  $G$ ,  $E$  and  $D$ 
2: repeat
3:   for number of training epochs do
4:     for number of mini-batches do
5:       // for discriminator  $D$ 
6:        $\theta_d \leftarrow \theta_d - \mu \frac{\partial \mathcal{L}_D(\mathbf{X}, \mathbf{I}; \theta_g, \theta_d)}{\partial \theta_d}$ 
7:       // for generator  $G$ 
8:        $\theta_g \leftarrow \theta_g - \mu \frac{\partial \mathcal{L}_G(\mathbf{X}, \mathbf{I}; \theta_g, \theta_d)}{\partial \theta_g}$ 
9:       // for extractor  $E$ 
10:       $\theta_e \leftarrow \theta_e - \mu \frac{\partial \mathcal{L}_{Rep}(\mathbf{I}, \mathbf{Y}; \theta_e)}{\partial \theta_e}$ 
11:      // for generator  $G$  and extractor  $E$ 
12:       $(\theta_g, \theta_e) \leftarrow (\theta_g, \theta_e) - \mu \frac{\partial \mathcal{L}_{Reg}(\mathbf{X}, \mathbf{Y}; \theta_g, \theta_e)}{\partial (\theta_g, \theta_e)}$ 
13:    end for
14:  end for
15: until convergence, got  $\hat{\theta}_g = \theta_g, \hat{\theta}_e = \theta_e, \hat{\theta}_d = \theta_d$ 
16: return  $\hat{\theta}_g, \hat{\theta}_e, \hat{\theta}_d$ 
```

Experiments

Networks	Abbreviation	Top1 Acc
MobileNetV2(T=6) [23]	Mb-6	92.84
MobileNetV2(T=10) [23]	Mb-10	93.94
MobileNetV2(T=16) [23]	Mb-16	94.29
ResNet-50 [9]	Res50	97.48
ResNet-101 [9]	Res101	98.12
DenseNet121 [11]	Dns121	97.45
DenseNet161 [11]	Dns161	97.70
PolyNetE [32]	Poly	98.46

Table 2: Identification results of different models on MegaFace dataset. ‘Top1 Acc’ refers to the top-1 face identification accuracy rate with 1M distractors.

Experiments

Architecture	Rec	Rec	Rep	Reg	Top1
	Adv	L2	L2	L2	Accuracy
FC	✗	✗	✗	✗	83.92
Arch1	✗	✓	✓	✗	85.65
Arch2	✓	✗	✓	✗	83.41
Arch3	✗	✗	✗	✓	94.05
Arch4	✗	✓	✗	✓	94.21
Arch5	✗	✓	✓	✓	94.93
Arch6	✓	✗	✗	✓	94.80
Arch7	✗	✗	✓	✓	NAN
R ³ AN	✓	✓	✓	✓	95.97

Table 3: Identification results of CMFR between MobileNetV2 (T=6) and ResNet-101 based on different architectures. Each row in this table is an architecture, and each column means a specific training process. The ‘✓’ and ‘✗’ represent for whether the architecture contains the process or not. ‘Rec:Adv’ and ‘Rec:L2’ means optimizing the generator by adversarial loss or L2 loss; ‘Rep:L2’ is extractor’s optimization; ‘Reg:L2’ represents the regression path. ‘Top1 Accuracy’ refers to the top-1 face identification accuracy rate with 1M distractors.

Source : MobileNet V2 → Target: ResNet-101

Experiments

Src	Tgt	Src	Tgt	FC Tgt	$R^3AN Tgt$
Mb-6	Mb-10	92.84	93.94	83.66	94.36
Mb-10	Mb-6	93.94	92.84	84.17	94.69
Mb-6	Mb-16	92.84	94.29	83.54	94.18
Mb-10	Mb-16	93.94	94.29	84.19	94.33
Mb-6	Res50	92.84	97.48	83.78	94.48
Mb-6	Res101	92.84	98.12	83.92	95.97
Mb-6	Poly	92.84	98.46	82.92	97.44
Mb-6	Dns161	92.84	97.70	82.95	95.66
Res50	Mb-6	97.48	92.84	87.49	97.60
Res101	Mb-6	98.12	92.84	88.17	98.19
Poly	Mb-6	98.46	92.84	88.51	98.34
Dns161	Mb-6	97.70	92.84	87.63	97.64
Res50	Res101	97.48	98.12	89.29	97.69
Res101	Res50	98.12	97.48	89.38	97.86
Dns121	Dns161	97.45	97.70	87.41	97.81
Res50	Poly	97.48	98.46	87.24	98.29

Table 5: Results of CMFR among different prior models. We use proposed models to map distribution of ‘Src’ (source model) to the distribution ‘Tgt’ (target model). The evaluation is established by taking the learned representation from the left model of ‘|’ as probe and right model’s output as gallery. ‘FC’ and ‘ R^3AN ’ refers to the architecture in the first and last row of Tab. 3. Results are the top-1 face identification accuracy rates with 1M distractors.

Visualization



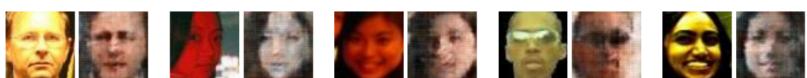
(a) Normal



(b) Same person



(c) Occlusion



(d) Lighting

Figure 4: The visualization of real faces and generated images from generator. Real faces are on the left, and generated images are on the right.

Visualization

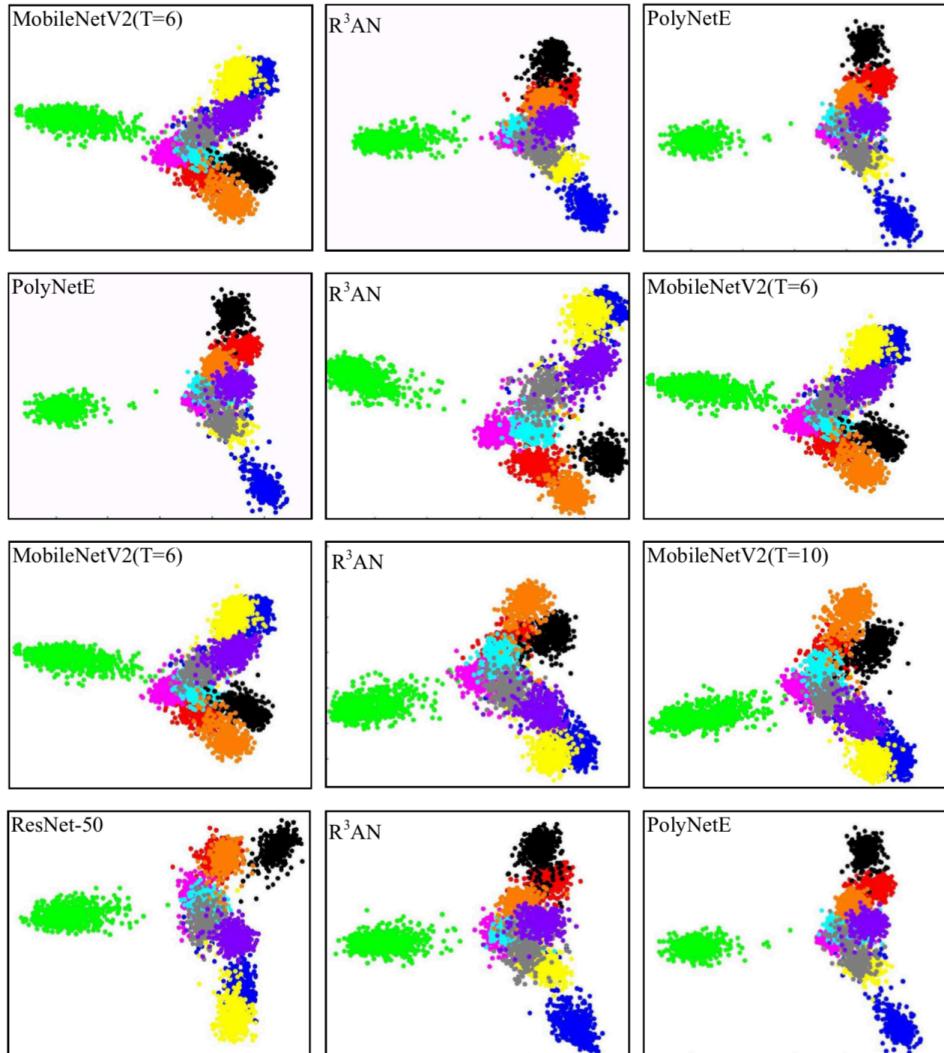


Figure 5: The visualization of feature distribution. From left to right in each row, images of distribution are from source model, R³AN and target model, respectively.

Thanks.

