



WWW.ECCV2020.EU



Black-Box Face Recovery from Identity Features

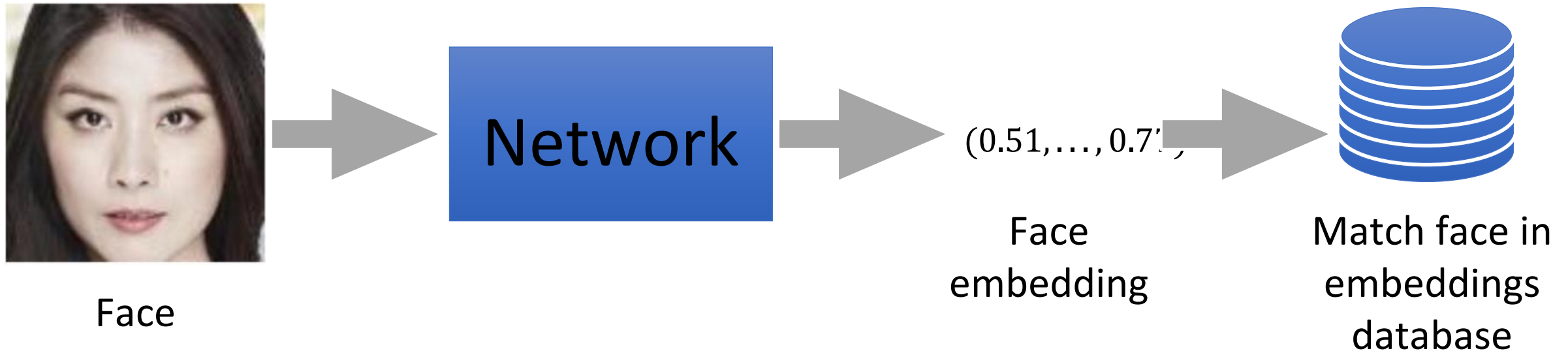
**Anton Razzhigaev^{1,2}, Klim Kireev^{1,2}, Edgar Kaziahmedov^{1,2}, Nurislam
Tursynbek^{1,2}, Aleksandr Petiushko^{1,3}**

Huawei¹, Skoltech², MSU³

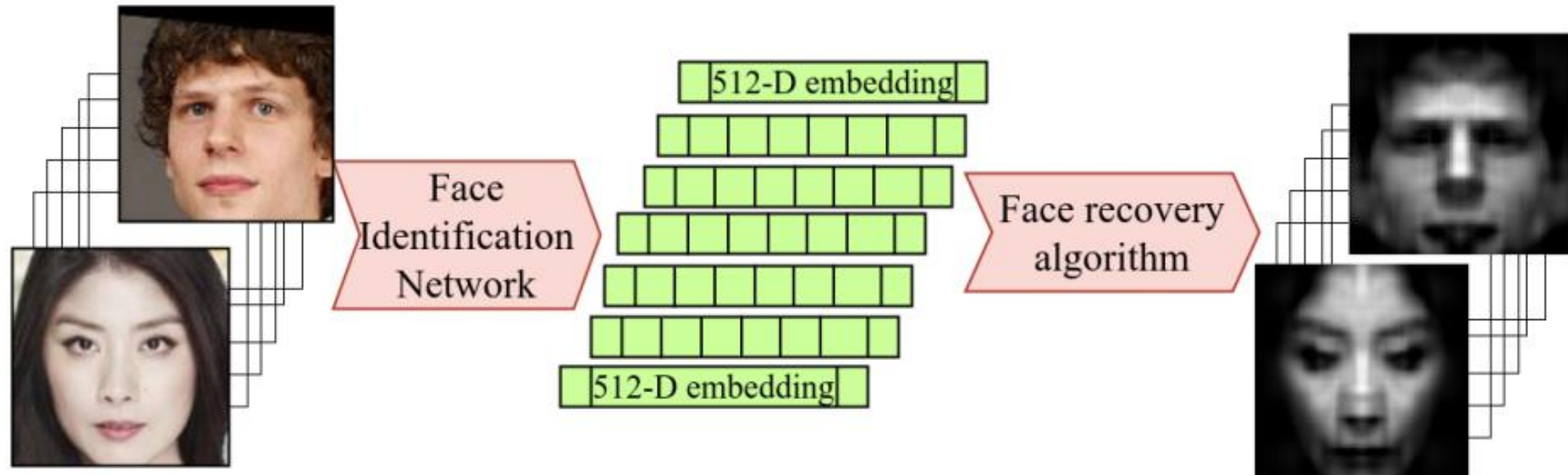


***Work done at the Huawei Moscow Research Center, Intelligent Systems & Data
Science lab**

Face recognition pipeline



Idea of face recovery



Novelty of our approach

Table 1: Comparison table

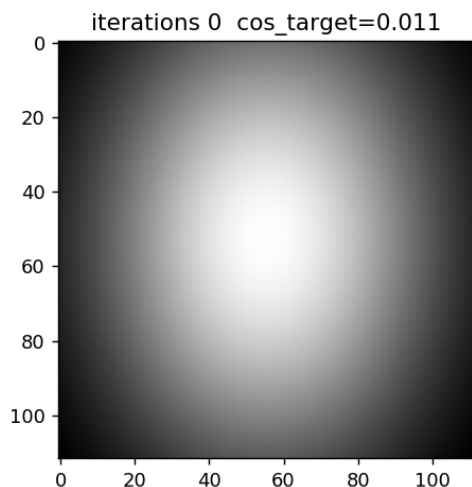
| Algorithm | Target model | Setting | Dataset-free |
|-------------------------|-------------------------------|-----------|--------------|
| Ours | Arcface output | Black-Box | + |
| NBNet[19] | FaceNet output | Black-Box | - |
| Cole et. al. [3] | FaceNet intermediate features | White-box | - |
| CNN[34] | FaceNet output features | White-box | - |
| Gradient wrt input [18] | Any classifier output | White-box | + |

- The **key point** of the proposed algorithm is that it **reconstructs not just faces but** recognizable **identities** without any prior knowledge about how faces look like
- It is not just face generator, it is a way to reveal appearance of a person from the face embedding

Algorithm



Original



Algorithm 1 Face recovery algorithm

INPUT: target face embedding y , black-box model M , loss function L , $N_{queries}$

```
1:  $X \leftarrow 0$ 
2: Initialize  $G_0$ 
3: for  $i \leftarrow 0$  to  $N_{queries}$  do:
4:   Allocate image batch  $\mathbf{X}$ 
5:   Sample batch  $\mathbf{G}$  of random gaussians
6:    $\mathbf{X}_j = X + G_0 + \mathbf{G}_j$ 
7:    $\mathbf{y}' = M(\mathbf{X})$ 
8:    $\text{ind} = \text{argmin} \left( L(\mathbf{y}'_i, y) \right)$ 
9:    $X \leftarrow X + \mathbf{G}_{\text{ind}}$ 
10:   $G_0 \leftarrow 0.99 \cdot G_0$ 
11:   $i \leftarrow i + \text{batchsize}$ 
12: end for
13:  $X \leftarrow X + G_0$ 
```

OUTPUT: reconstructed face X

$$L(y, y') = \lambda \cdot (\|y\| - \|y'\|)^2 - s(y, y'),$$

where,

s - cosine similarity function,

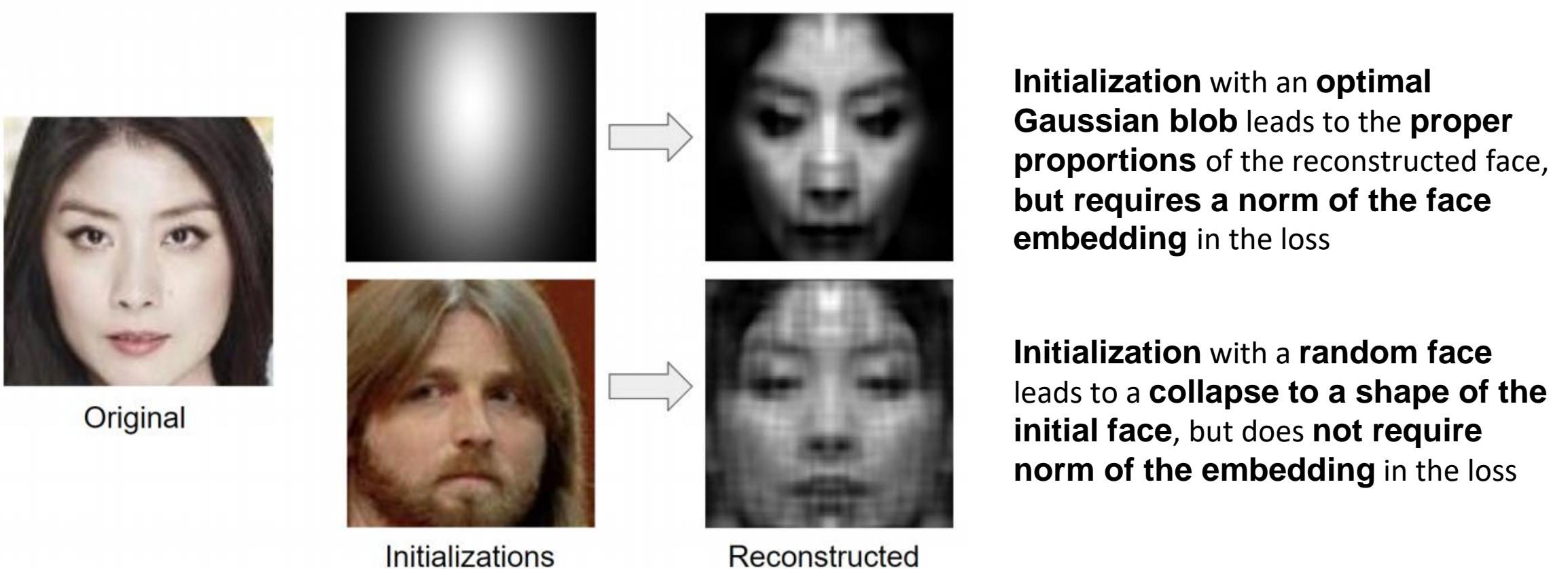
$\|y\|$ - L₂ norm of the target embedding,

$\|y'\|$ - L₂ norm of the embedding of

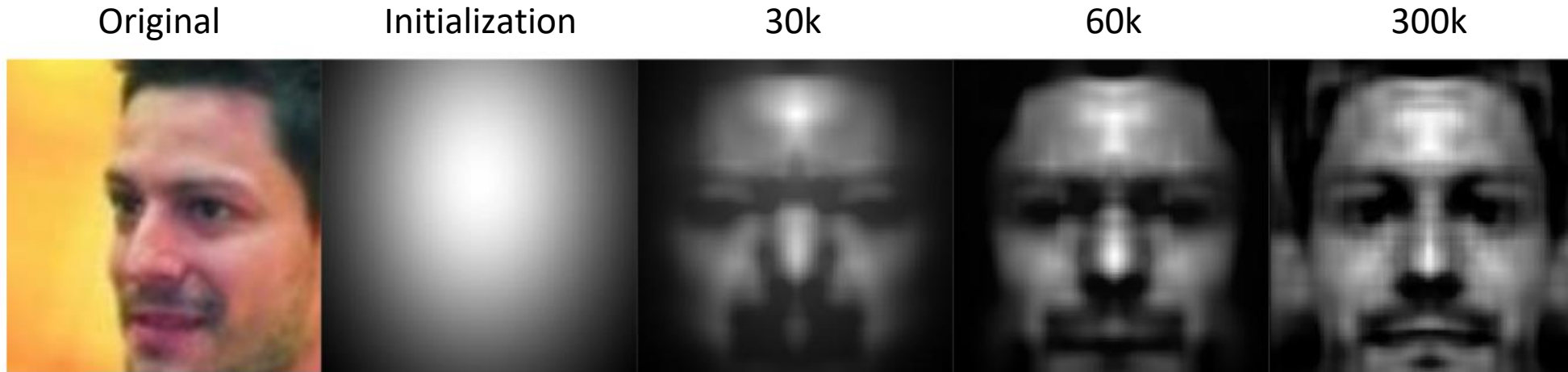
a reconstructed image,

$\lambda = 0.0025$, empirically found hyperparameter

Different initialization techniques



Function space for “drawing”



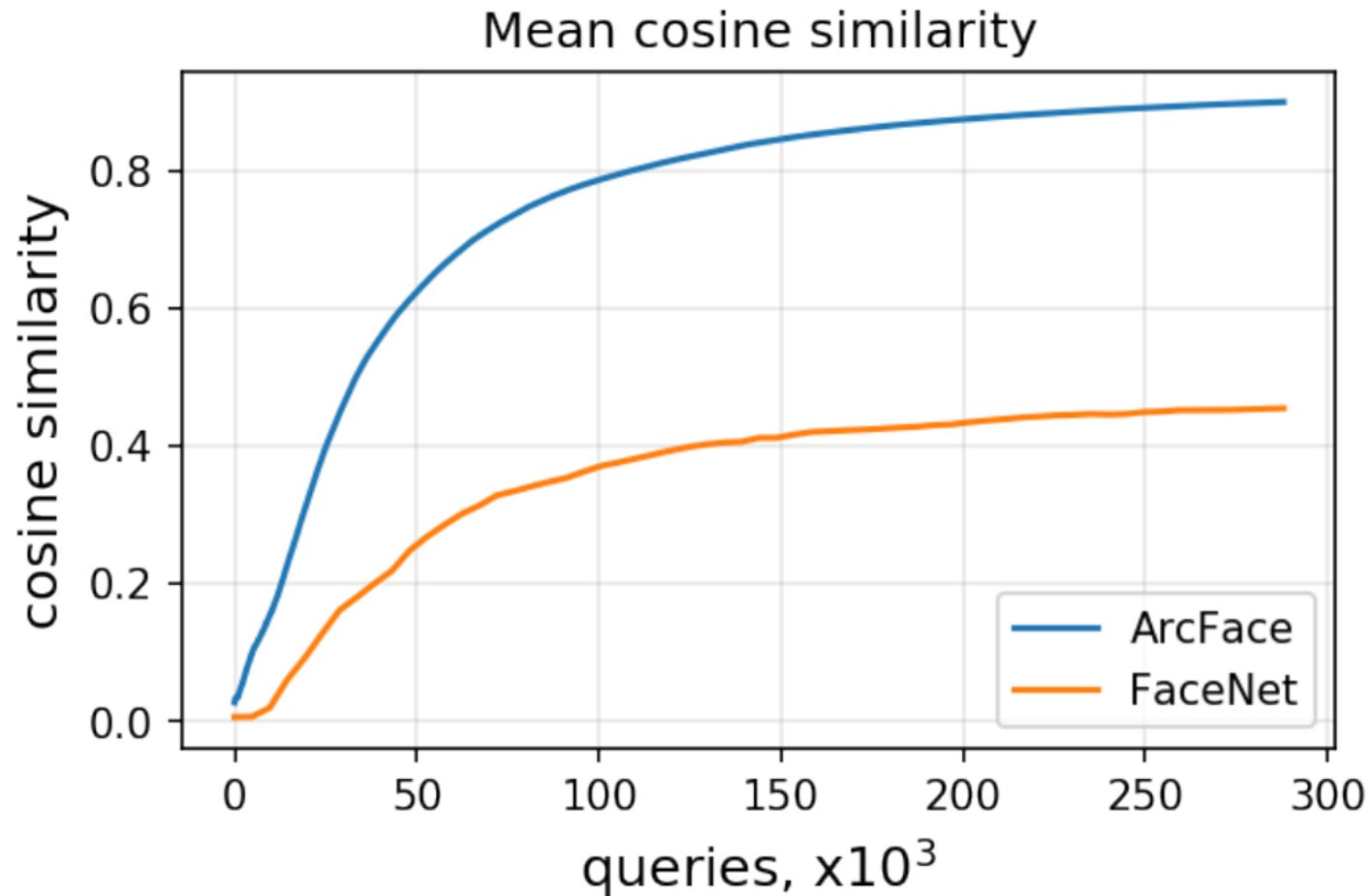
$$G(x, y) = A \cdot \exp \frac{(x - x_0)^2}{2\sigma_1^2} \exp \frac{(y - y_0)^2}{2\sigma_2^2}$$

x, y - pixel coordinates in the image,
 x_0, y_0 - coordinates of a center of gaussian,
 σ_1, σ_2 - vertical and horizontal standard deviations,
 A - amplitude

Gaussian 2D functions:

1. Semi-local
2. Frequency control through σ

Evaluation



- Attacking: **ArcFace** network
- Evaluating: using an “**independent critic**” – **FaceNet** network

Symmetry VS asymmetry VS color

Symmetry, black&white

Asymmetry, black&white

Symmetry, rgb



Original

ArcFace: 0.978

FaceNet: 0.721

ArcFace: 0.992

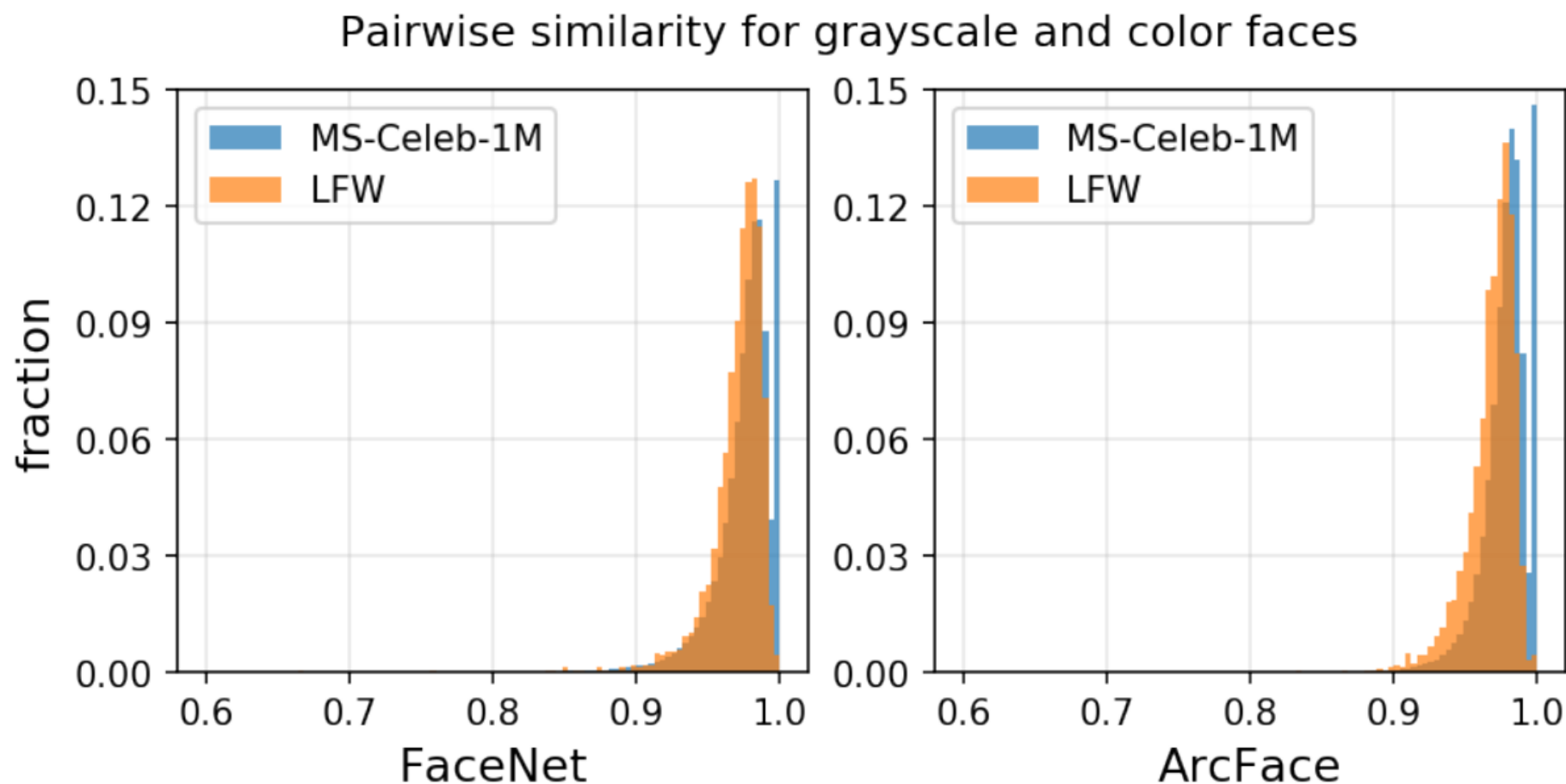
FaceNet: 0.685

ArcFace: 0.961























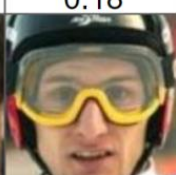

FaceNet: 0.314

Best results with **symmetry** and **black & white** constraints

Identity features do not contain much information about color



Comparison with others

| | | | | | | |
|--------------|---|---|---|---|--|---|
| Our method: |  |  |  |  |  |  |
| ArcFace: | 0.97 | 0.97 | 0.94 | 0.97 | 0.85 | 0.73 |
| FaceNet: | 0.70 | 0.75 | 0.72 | 0.78 | 0.38 | -0.09 |
| NBNet (WB): |  |  |  |  |  |  |
| ArcFace: | 0.17 | 0.21 | 0.12 | 0.26 | 0.06 | 0.09 |
| FaceNet: | 0.02 | 0.32 | 0.25 | 0.46 | -0.01 | 0.35 |
| NBNet (RGB): |  |  |  |  |  |  |
| ArcFace: | 0.28 | 0.46 | 0.34 | 0.54 | 0.12 | 0.21 |
| FaceNet: | 0.59 | 0.53 | 0.44 | 0.74 | 0.18 | 0.41 |
| Original: |  |  |  |  |  |  |

Our method:

- produces **less realistic faces**,
- but more likely **preserves identity of a person** which **correlates** with similarities of “**independent critic**” network

NBNet: Mai, G., Cao, K., Yuen, P.C., Jain, A.K.: *On the reconstruction of face images from deep face templates*. IEEE transactions on pattern analysis and machine intelligence, 2018

Comparison with others

| Method | ArcFace | FaceNet | # of queries |
|---|-------------|-------------|--------------|
| (Ours) Symmetric gauss, LFW (wb) | 0.92 | 0.46 | 300k |
| (Ours) Asymmetric gauss, LFW (wb) | 0.85 | 0.42 | 400k |
| NBNet, LFW (RGB) | 0.25 | 0.34 | 3M |
| NBNet, LFW (wb) | 0.19 | 0.27 | 3M |
| (Ours) Symmetric gauss, MS1M-ArcFace (wb) | 0.91 | 0.44 | 300k |
| NBNet, MS1M-ArcFace (RGB) | 0.26 | 0.38 | 3M |
| NBNet, MS1M-ArcFace (wb) | 0.20 | 0.32 | 3M |

Table 2: Average cosine similarity by ArcFace and FaceNet (independent critic) between embedding of a reconstructed image and embedding of target image for subsets of 1000 images from LFW and MS1M-ArcFace and corresponding number of queries.

Conclusions

- We demonstrate that it is **possible** to **recover recognizable faces** from deep **feature vectors** of a face-recognition model in a **black-box mode** with **no prior knowledge**.
- The **proposed method outperforms** current solutions in terms of the **average cosine similarity** of embeddings produced by the **attacked model** and an **independent critic**.
- The **proposed method** requires a significantly **smaller number of queries** compared to previous solutions and **does not need** prior information such as **proper training dataset**.