# Face Recognition at a Distance

## Prof. Xiaoming Liu

Computer Vision Lab
Department of Computer Science and Engineering
Michigan State University

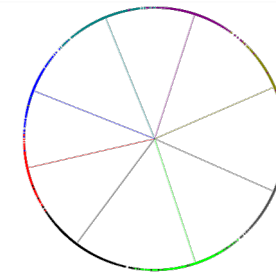# Tremendous Research Progress
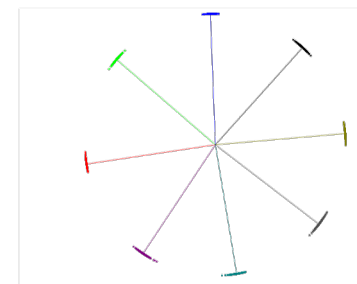

DR-GAN


Age synthesis

Age: 27


Disguised faces


Softmax
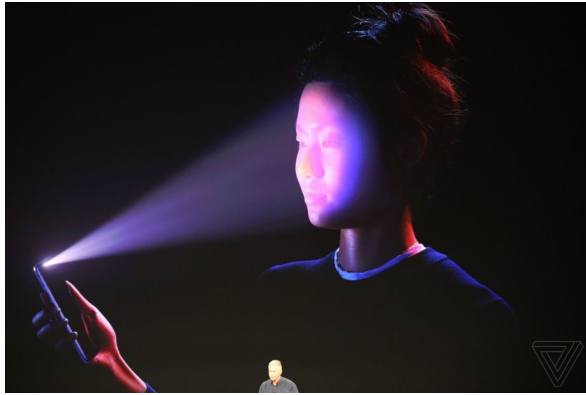
SphereFace

CosFace


ArcFace

L. Tran et. al., Representation Learning by Rotating Your Faces, PAMI, 2018
H. Yang et. al., Learning Continuous Face Age Progression: A Pyramid of GANs, PAMI, 2019.
M. Singh et. al., Recognizing disguised faces in the wild. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 2019.
J. Deng et. al., ArcFace: Additive Angular Margin Loss for Deep Face Recognition, CVPR, 2019.

# Successful Applications


Apple


Alipay


Boarding in Airports


Entrance to Beijing University

**Objective**: Recognize individuals from a video stream captured at a distance and altitude.

**Modality**: Face, gait and body



Biometrics

**Outline**:

➢ Generic matcher: AdaFace (CVPR 2022)

➢ Domain adaption: CFSM (ECCV 2022)

➢ Video-based recognition: CAFace (NeurIPS 2022)
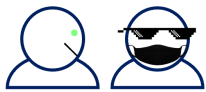
➢ 3D body matching (Under review)

# AdaFace: Quality Adaptive Margin for Face Recognition

Minchul Kim, Anil K. Jain, Xiaoming Liu
CVPR 2022

# Problem Definition

## Training Datasets have Varying **Qualities**
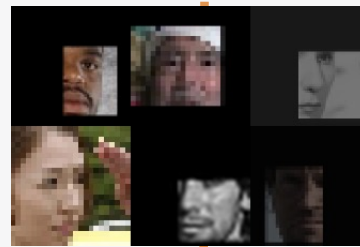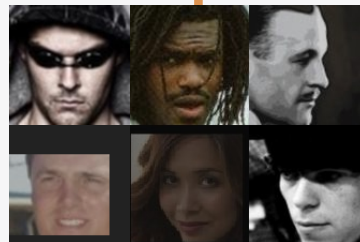
### Pose and Occlusion

Faces that are front facing and free of occlusions such as hands or sunglasses are identifiable.

### Landmark & Key points

Images with visible and detectable facial landmarks are identifiable.

Easy to Recognize

Hard to Recognize

### Blur

Subject's distance, camera setting and other environmental factors cause the image to be blurred.

**H** **W**

### Size

Too low image resolution causes the subject to be unidentifiable.

### Illumination

Too dark or too bright images cause the subject to be unidentifiable.

**Source of Problem (Impossible to recognize)**

Training dataset **without identifiable traits** can be equivalent to **noisy label samples**

**Easy to Recognize**

**Hard to Recognize**

## Adaptive Sample Emphasis During Training

**Easy Samples** — Easy samples are well classified

**Difficult Samples** — Hard sample mining

**Unidentifiable Samples** — Avoid learning from bad samples

Low ➝ High

**Magnitude of Emphasis**

# One More Way to Look at an Image

High Quality

Quality

Low Quality

Hypothesis:
**Difficulty** distribution will be different based on image quality.



Easy to Recognize ......... Hard to Recognize ......... Impossible to Recognize

**Difficulty**

# Our Findings and Methods



- Emphasis Magnitude
- High Quality
- Quality
- Low Quality

High Emphasis

Hard samples are emphasized

**2** **Feature norm** can approximate the quality.

**1** **Margin Functions** can change the slope.

High Emphasis

Unrecognizable samples are avoided

Low Emphasis

Difficulty

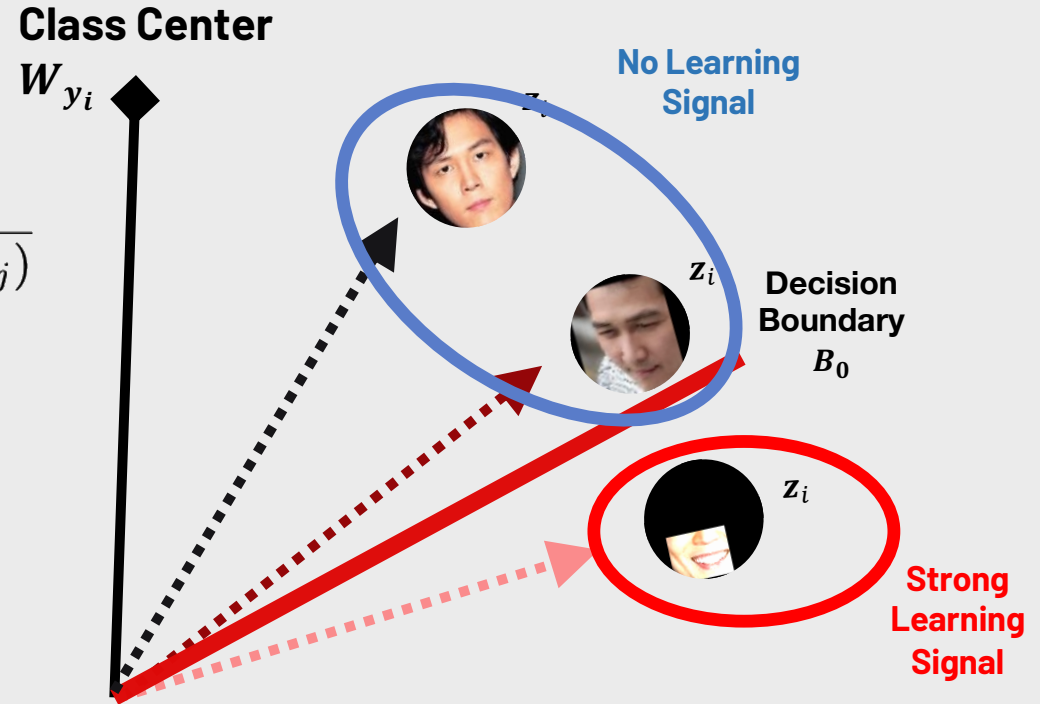Easy to Recognize        Hard to Recognize        Impossible to Recognize

# Effect of Margin on Sample Emphasis



**Margin-based SoftMax Loss**

$$\mathcal{L} = -\log \frac{\exp(f(\theta_{y_i}, m))}{\exp(f(\theta_{y_i}, m)) + \sum_{j \neq y_i}^{n} \exp(s \cos \theta_j)}$$
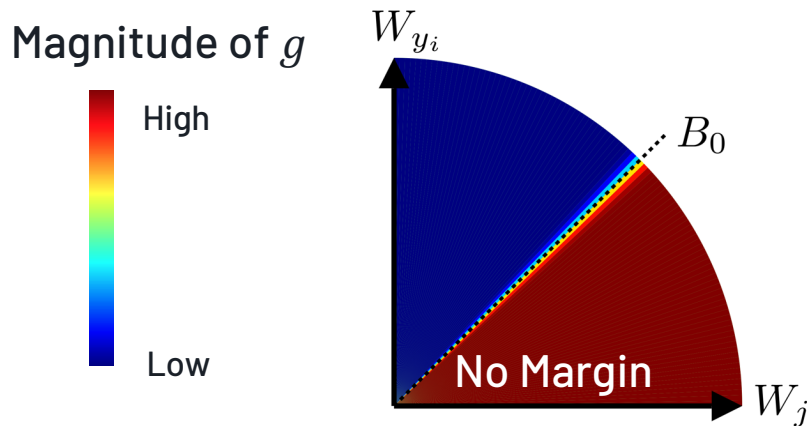
$$f(\cdot)_{\text{Additive}} = \begin{cases} s((\cos \theta_{y_i}) - m) & j = y_i \\ s \cos \theta_{y_i} & j \neq y_i \end{cases}$$

$$f(\cdot)_{\text{Angular}} = \begin{cases} s \cos(\theta_{y_i} + m) & j = y_i \\ s \cos \theta_{y_i} & j \neq y_i \end{cases}$$

**Class Center** $W_{y_i}$

$z_i$ — **No Learning Signal**

$z_i$ — **Decision Boundary** $B_0$

$z_i$ — **Strong Learning Signal**

**Plot of Gradient Scaling Term**

$$\frac{\partial \mathcal{L}_{\text{CE}}}{\partial \boldsymbol{x}_i} = \sum_{k=1}^{C} \left( P_k^{(i)} - \mathbb{1}(y_i = k) \right) \frac{\partial f(\cos \theta_k)}{\partial \cos \theta_k} \frac{\partial \cos \theta_k}{\partial \boldsymbol{x}_i}.$$

Magnitude of $g$

High

Low
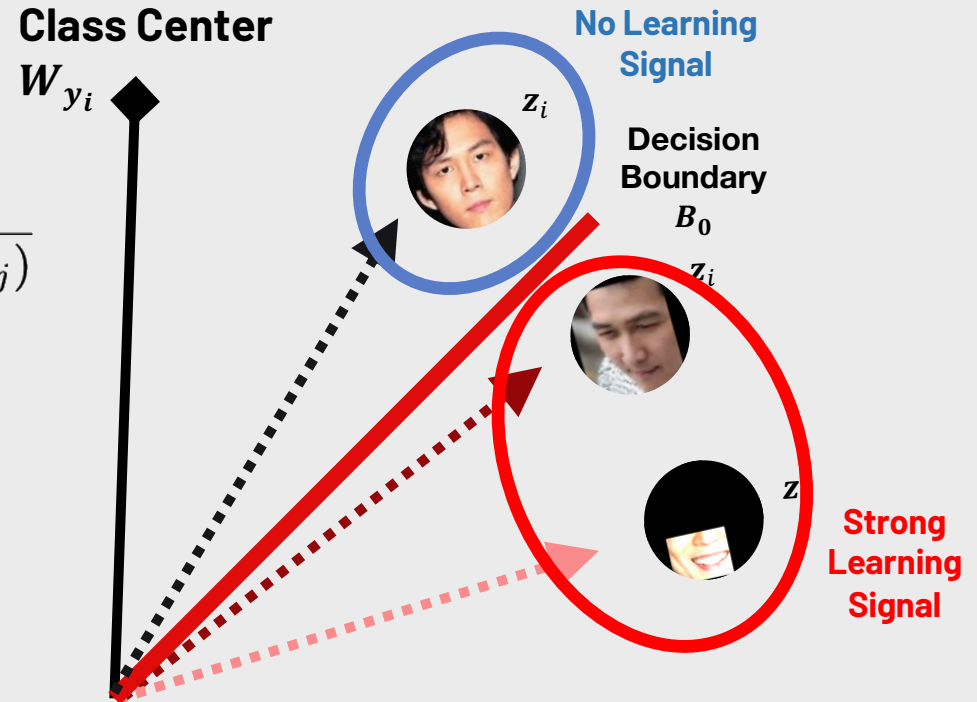
$W_{y_i}$

$B_0$

No Margin

$W_j$

# Effect of Margin on Sample Emphasis

**Margin-based SoftMax Loss**

$$\mathcal{L} = -\log \frac{\exp(f(\theta_{y_i}, m))}{\exp(f(\theta_{y_i}, m)) + \sum_{j \neq y_i}^{n} \exp(s\cos\theta_j)}$$
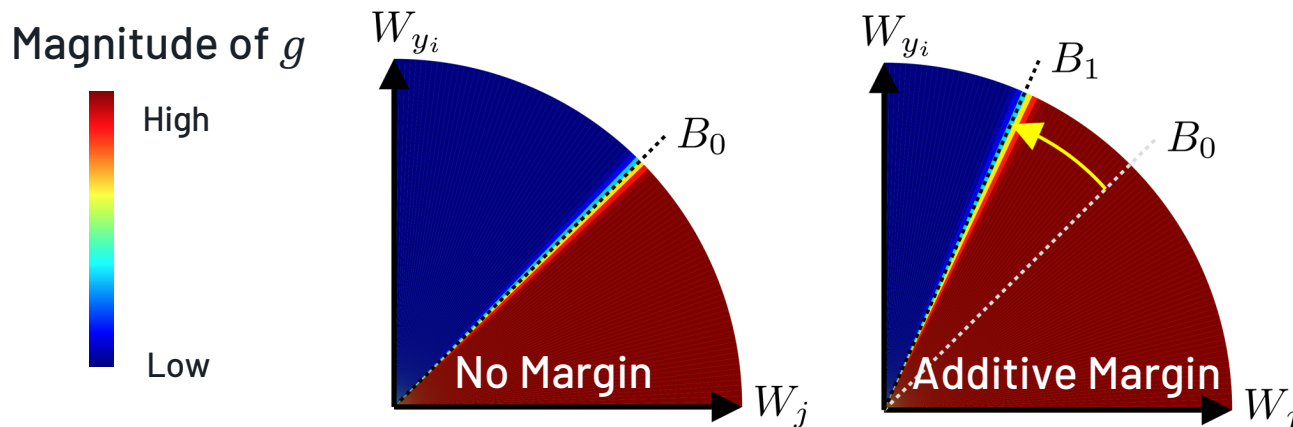
$$f(\cdot)_{\text{Additive}} = \begin{cases} s((\cos\theta_{y_i}) - m) & j = y_i \\ s\cos\theta_{y_i} & j \neq y_i \end{cases}$$

$$f(\cdot)_{\text{Angular}} = \begin{cases} s\cos(\theta_{y_i} + m) & j = y_i \\ s\cos\theta_{y_i} & j \neq y_i \end{cases}$$

**Class Center**

$W_{y_i}$

**No Learning Signal**

$z_i$

**Decision Boundary** $B_0$

$z_i$

$z$

**Strong Learning Signal**

**Plot of Gradient Scaling Term**

$$\frac{\partial \mathcal{L}_{\text{CE}}}{\partial \boldsymbol{x}_i} = \sum_{k=1}^{C} \Big(P_k^{(i)} - \mathbb{1}(y_i = k)\Big) \frac{\partial f(\cos\theta_k)}{\partial \cos\theta_k} \frac{\partial \cos\theta_k}{\partial \boldsymbol{x}_i}.$$

Magnitude of $g$

High

Low

$W_{y_i}$

$B_0$

No Margin

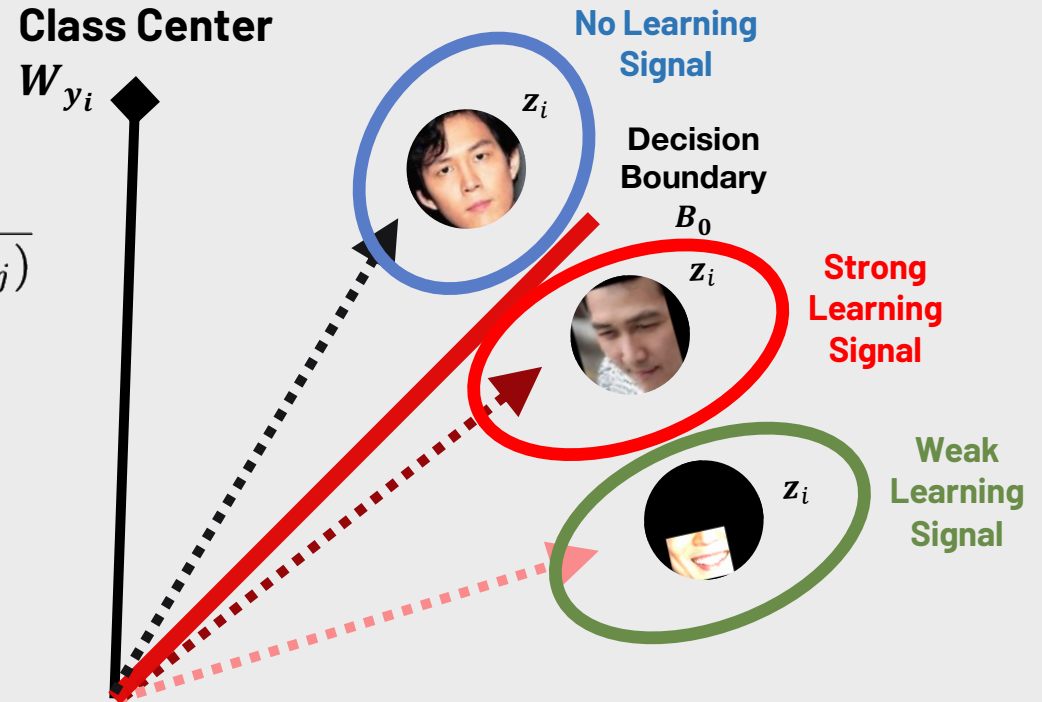$W_j$

$W_{y_i}$

$B_1$

$B_0$

Additive Margin

$W_j$

# Effect of Margin on Sample Emphasis

**Margin-based SoftMax Loss**

$$\mathcal{L} = -\log \frac{\exp(f(\theta_{y_i}, m))}{\exp(f(\theta_{y_i}, m)) + \sum_{j \neq y_i}^{n} \exp(s \cos \theta_j)}$$
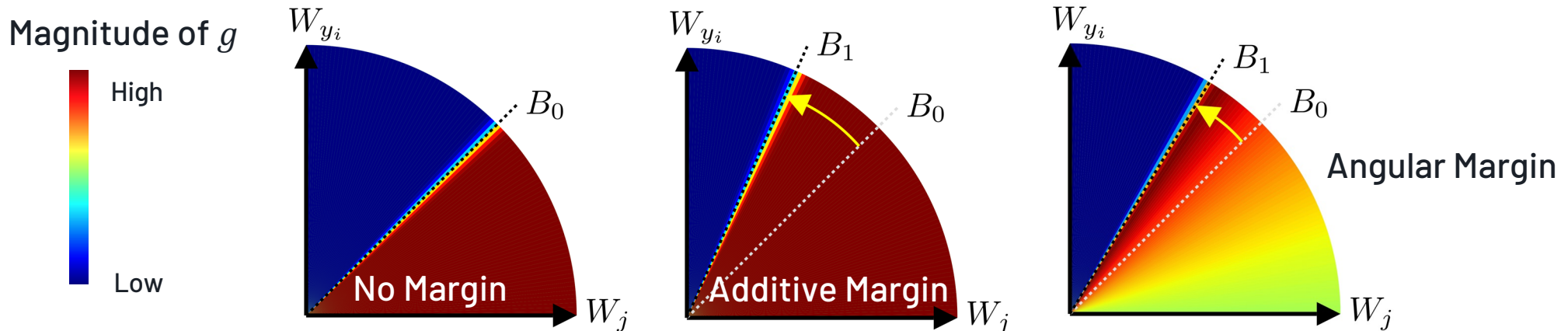
$$f(\cdot)_{\text{Additive}} = \begin{cases} s((\cos \theta_{y_i}) - m) & j = y_i \\ s \cos \theta_{y_i} & j \neq y_i \end{cases}$$

$$f(\cdot)_{\text{Angular}} = \begin{cases} s \cos(\theta_{y_i} + m) & j = y_i \\ s \cos \theta_{y_i} & j \neq y_i \end{cases}$$
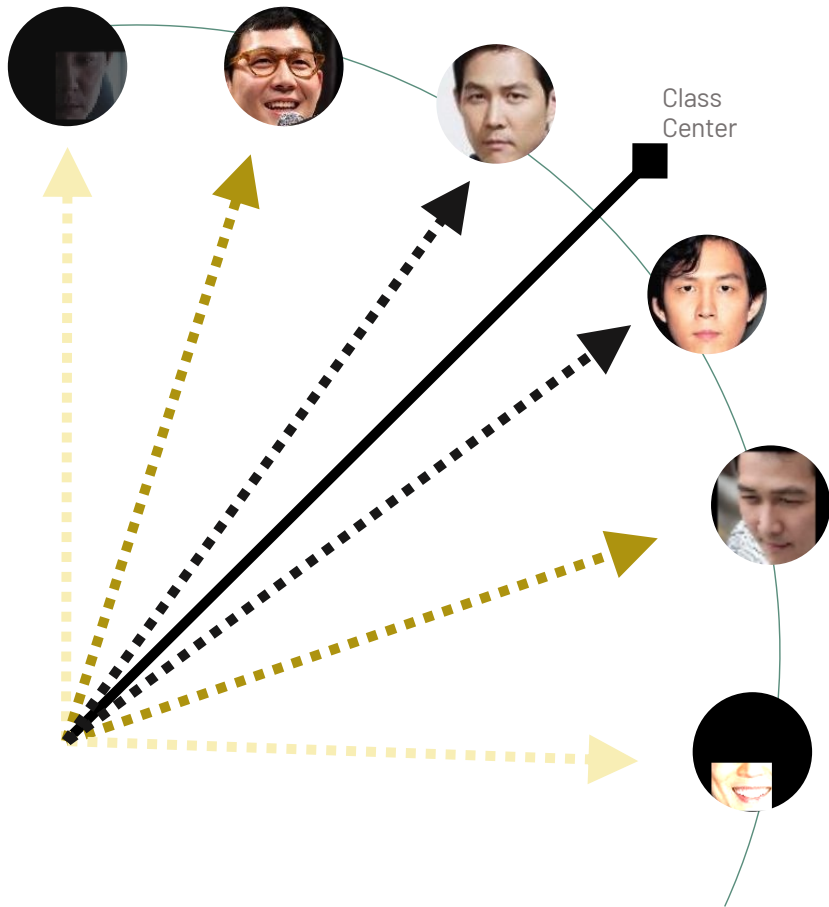
**Class Center**
$W_{y_i}$

**No Learning Signal**
$z_i$

**Decision Boundary**
$B_0$

**Strong Learning Signal**
$z_i$

**Weak Learning Signal**
$z_i$

**Plot of Gradient Scaling Term**

$$\frac{\partial \mathcal{L}_{\text{CE}}}{\partial \boldsymbol{x}_i} = \sum_{k=1}^{C} \left( P_k^{(i)} - \mathbb{1}(y_i = k) \right) \frac{\partial f(\cos \theta_k)}{\partial \cos \theta_k} \frac{\partial \cos \theta_k}{\partial \boldsymbol{x}_i}.$$
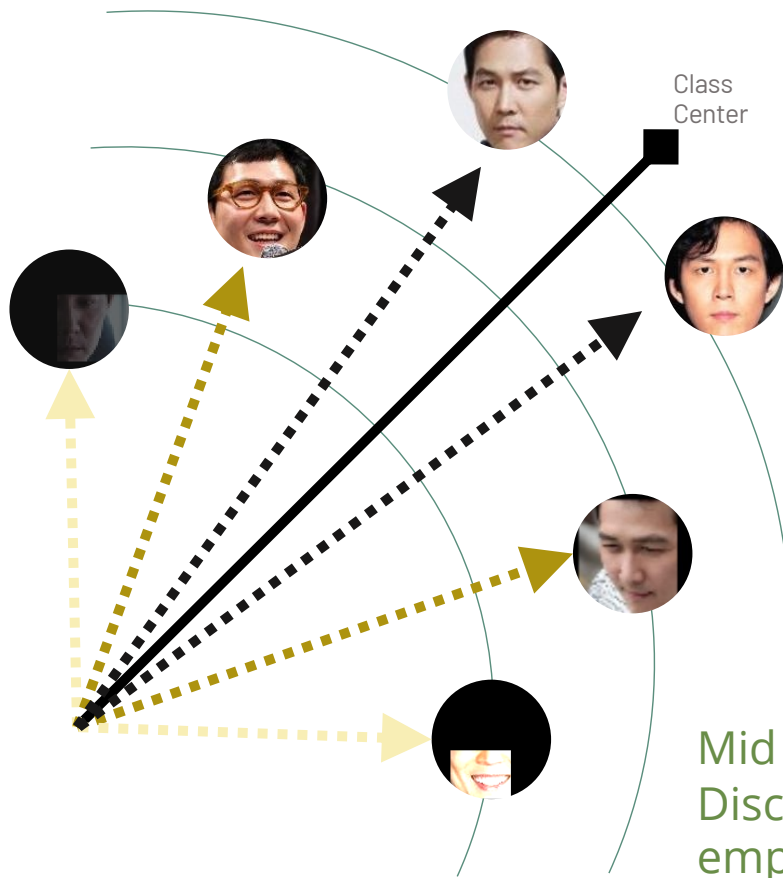
Magnitude of $g$

High

Low

$W_{y_i}$    $B_0$    No Margin    $W_j$

$W_{y_i}$    $B_1$    $B_0$    Additive Margin    $W_j$

$W_{y_i}$    $B_1$    $B_0$    Angular Margin    $W_j$

Feature Space

Class Center

Previous works apply
same margin for all samples

Unit Sphere Representation

# Method



Class Center

AdaFace
Adaptive Margin

High Norm = Negative Angular Margin
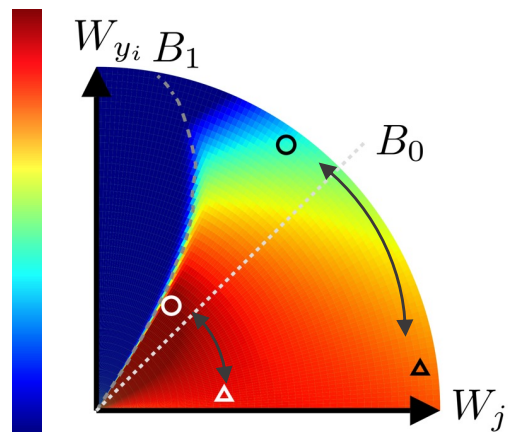De-emphasize trivial samples

Mid Norm = Additive Margin
Discriminative feature, equal
emphasis.

Low Norm = Positive Angular Margin
De-emphasize unrecognizable images

How do we emphasize different samples?

## AdaFace Objective

**Magnitude of g**



$$\mathcal{L} = -\log \frac{\exp(f(\theta_{y_i}, m))}{\exp(f(\theta_{y_i}, m)) + \sum_{j \neq y_i}^{n} \exp(s \cos \theta_j)}$$

$$f(\theta_j, m)_{\text{AdaFace}} = \begin{cases} s \cos(\theta_j + g_{\text{angle}}) - g_{\text{add}} & j = y_i \\ s \cos \theta_j & j \neq y_i \end{cases}$$
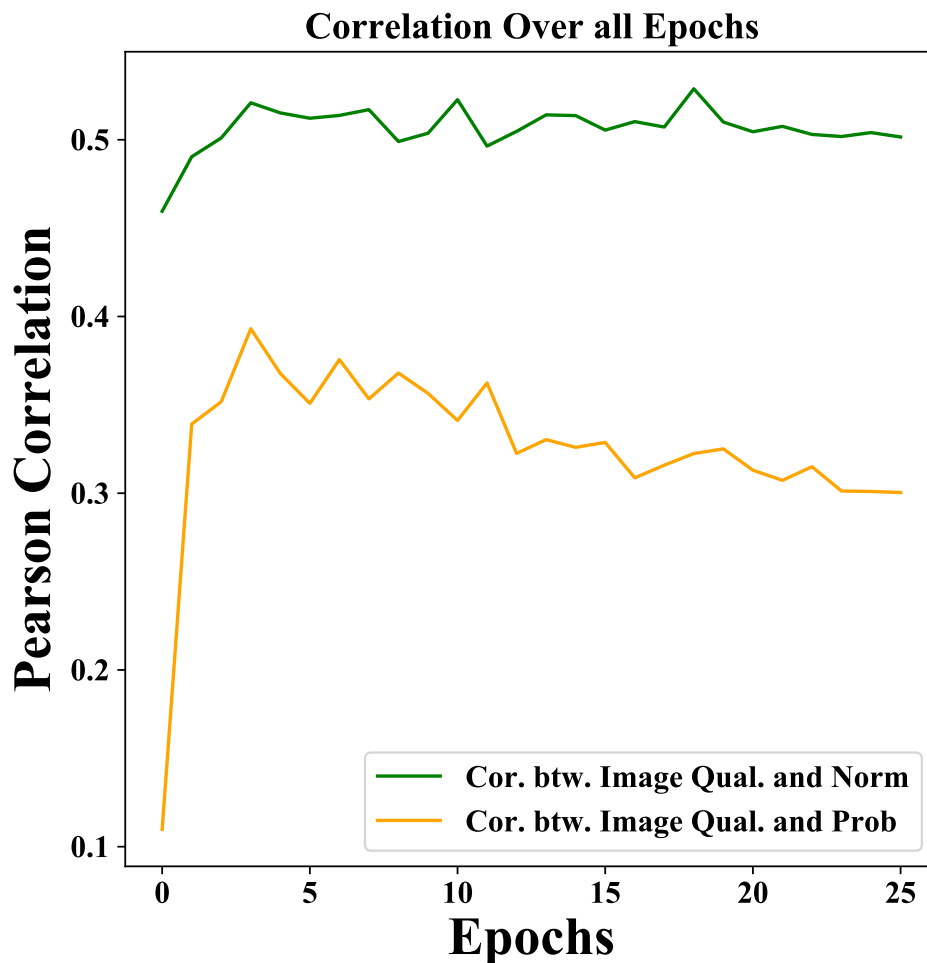
$$g_{\text{angle}} = -m \cdot \widehat{\|\boldsymbol{z}_i\|}, \quad g_{\text{add}} = m \cdot \widehat{\|\boldsymbol{z}_i\|} + m.$$

| ○ | Easier Sample, Low Norm | ○ | Easier Sample, High Norm |
|---|---|---|---|
| △ | Harder Sample, Low Norm | △ | Harder Sample, High Norm |

$$\widehat{\|\boldsymbol{z}_i\|} = \left\lfloor \frac{\|\boldsymbol{z}_i\| - \mu_z}{\sigma_z/h} \right\rfloor_{-1}^{1}$$

Combine different margin functions adaptively
to emphasize samples of different difficulty based on the image quality.
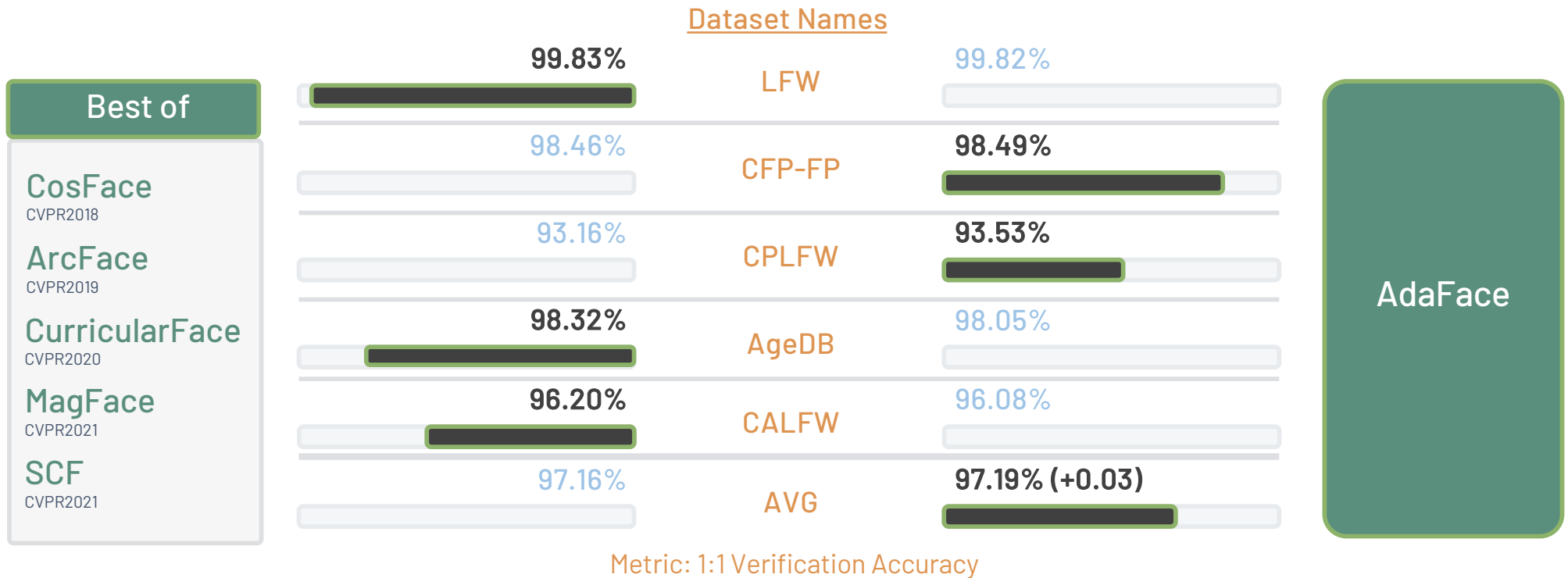
# Relationship between IQ and Feature Norm



Correlation Over all Epochs

- Image Quality: calculated with BRISQUE algorithm.

- The correlation between feature norm and the image quality exists from the early stage of training

**We use feature norm as a proxy for the quality and use it to change the margin.**
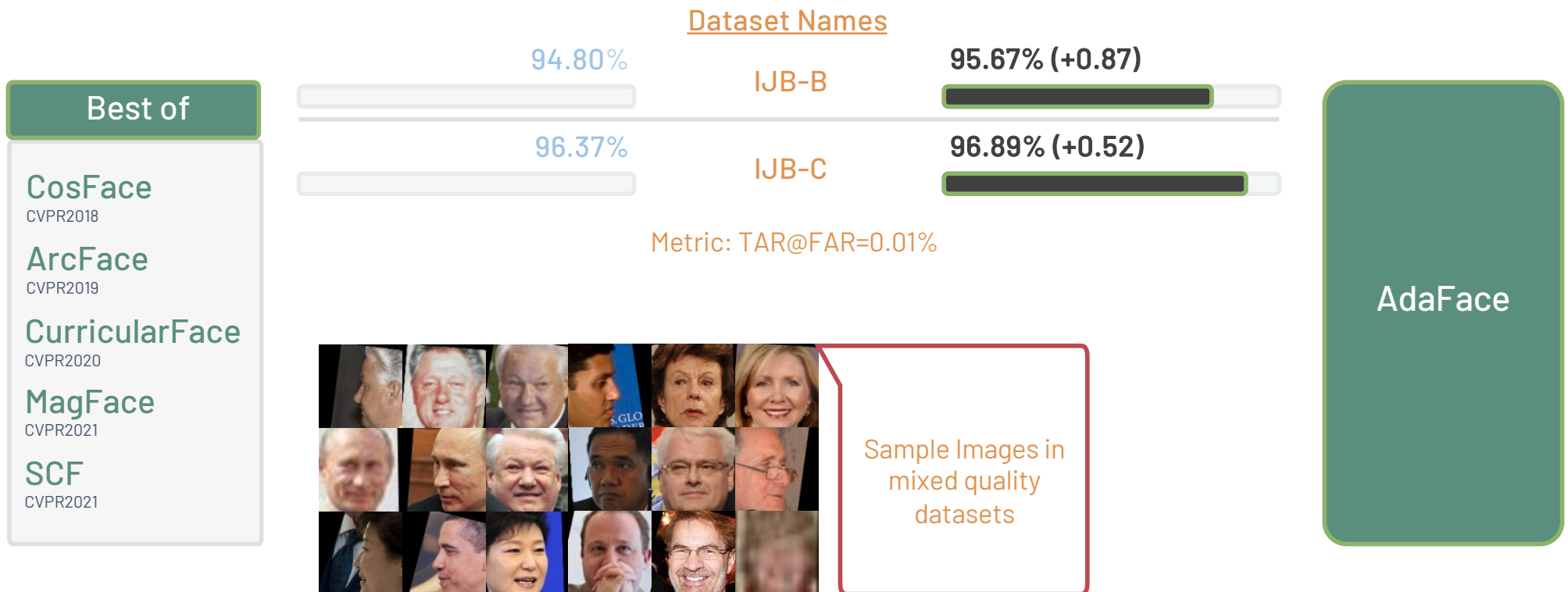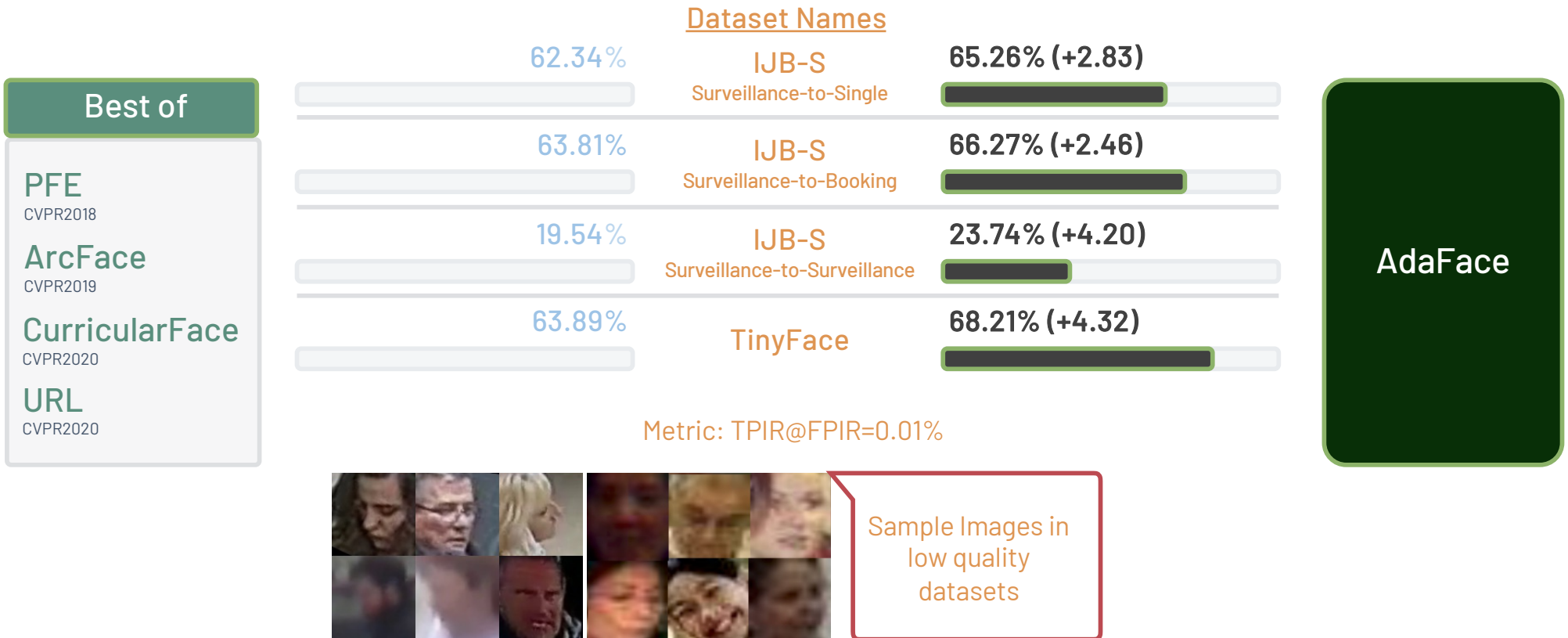
# Performance in High Quality Datasets

**Dataset Names**

| Best of | | AdaFace |
|---|---|---|
| **CosFace** CVPR2018 | LFW — 99.83% | 99.82% |
| **ArcFace** CVPR2019 | CFP-FP — 98.46% | **98.49%** |
| **CurricularFace** CVPR2020 | CPLFW — 93.16% | **93.53%** |
| **MagFace** CVPR2021 | AgeDB — **98.32%** | 98.05% |
| **SCF** CVPR2021 | CALFW — **96.20%** | 96.08% |
| | AVG — 97.16% | **97.19% (+0.03)** |

Metric: 1:1 Verification Accuracy



Sample Images in high quality datasets

# Performance in Mixed Quality Datasets



Best of

CosFace
CVPR2018

ArcFace
CVPR2019

CurricularFace
CVPR2020

MagFace
CVPR2021

SCF
CVPR2021

Dataset Names

94.80%   IJB-B   95.67% (+0.87)

96.37%   IJB-C   96.89% (+0.52)

Metric: TAR@FAR=0.01%

Sample Images in mixed quality datasets

AdaFace

# Performance in Low Quality Datasets

| Best of | | Dataset Names | | AdaFace |
|---|---|---|---|---|
| | 62.34% | **IJB-S**<br>Surveillance-to-Single | **65.26% (+2.83)** | |
| **PFE**<br>CVPR2018 | 63.81% | **IJB-S**<br>Surveillance-to-Booking | **66.27% (+2.46)** | |
| **ArcFace**<br>CVPR2019 | 19.54% | **IJB-S**<br>Surveillance-to-Surveillance | **23.74% (+4.20)** | |
| **CurricularFace**<br>CVPR2020 | 63.89% | **TinyFace** | **68.21% (+4.32)** | |
| **URL**<br>CVPR2020 | | | | |

Metric: TPIR@FPIR=0.01%



Sample Images in low quality datasets

# Controllable and Guided Face Synthesis for Unconstrained Face Recognition

Feng Liu, Minchul Kim, Anil Jain, and Xiaoming Liu
ECCV 2022

# Unconstrained Face Recognition

➤ Domain gap between the semi-constrained training datasets and unconstrained testing scenarios.



*Large-scale* Training datasets
VGG2. WebFace

Testing scenarios
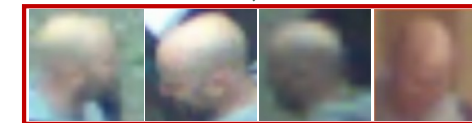LFW                    IJB-S

Semi-constrained
Faces collected from the web.

Semi-constrained
Accuracy:99.83%

Unconstrained
Accuracy<70%

➤ Potential solution

*Source domain*                                        *Target unconstrained domain*



Face Synthesis    →    False    Discriminator    True
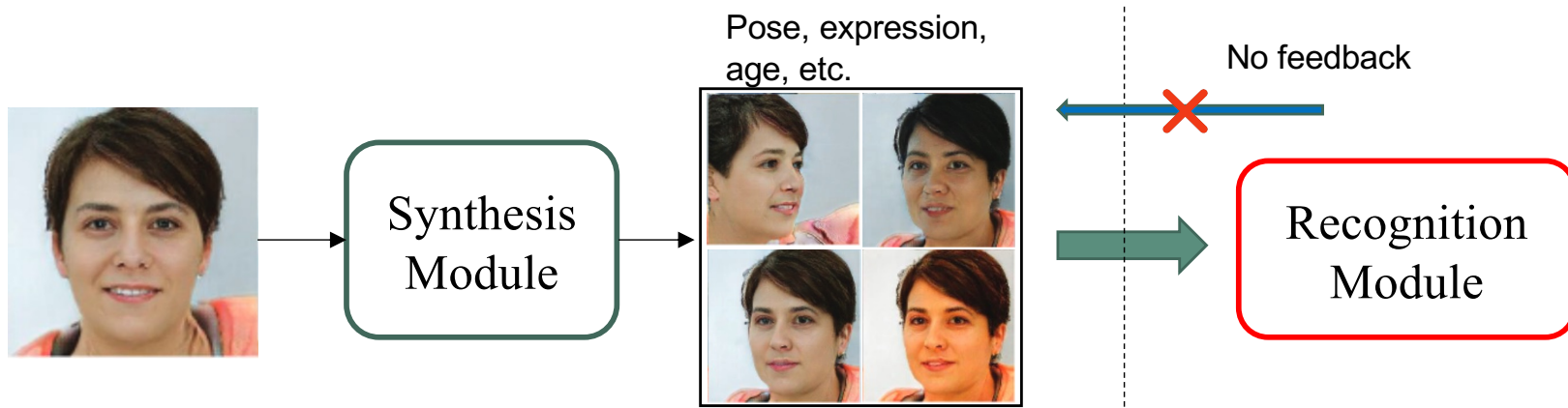
- Low resolution
- Motion blurring
- Bad illumination
- Turbulence effect

......

# Motivation

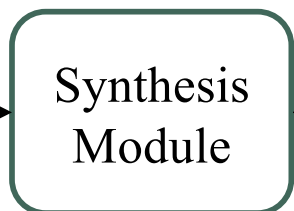➢ Previous synthesis models: ***limited face properties***; offline and ***blind*** data augmentation.
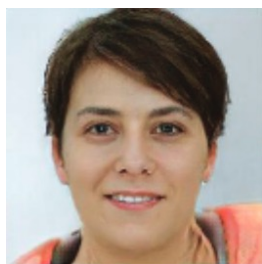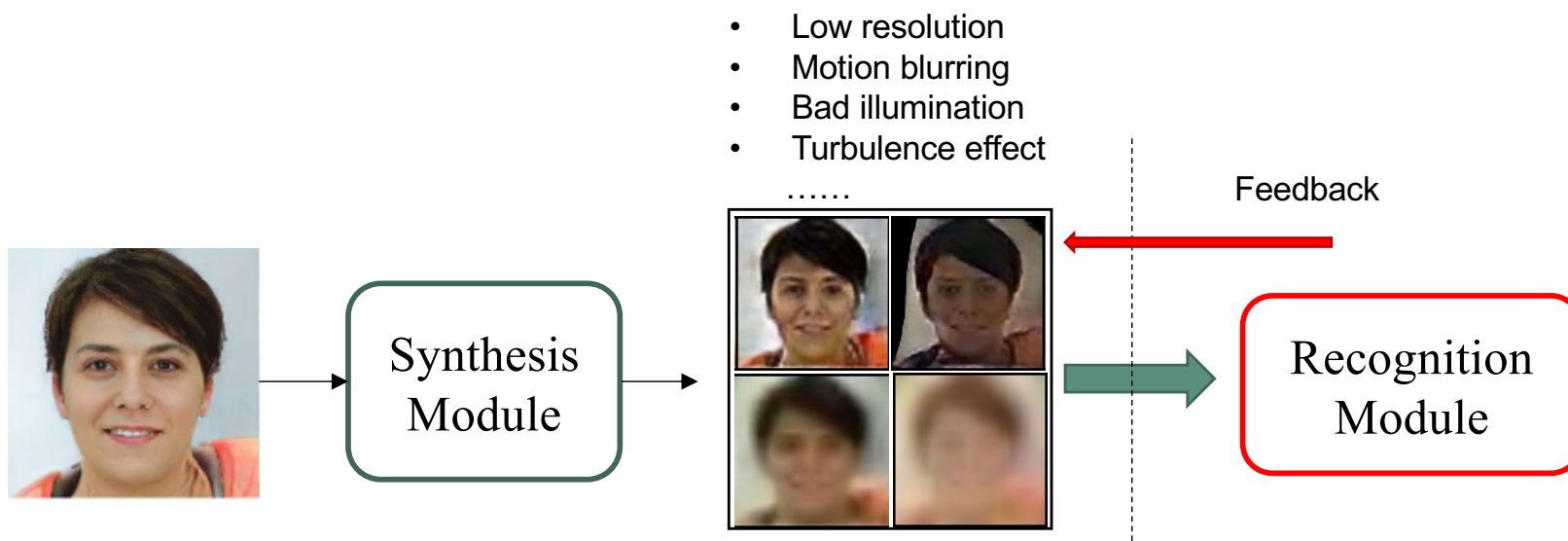
# Motivation

➢ Previous synthesis models: ***limited face properties***; offline and ***blind*** data augmentation.

➢ Facial properties should be generalizable to the challenging unconstrained testing scenarios.

- Low resolution
- Motion blurring
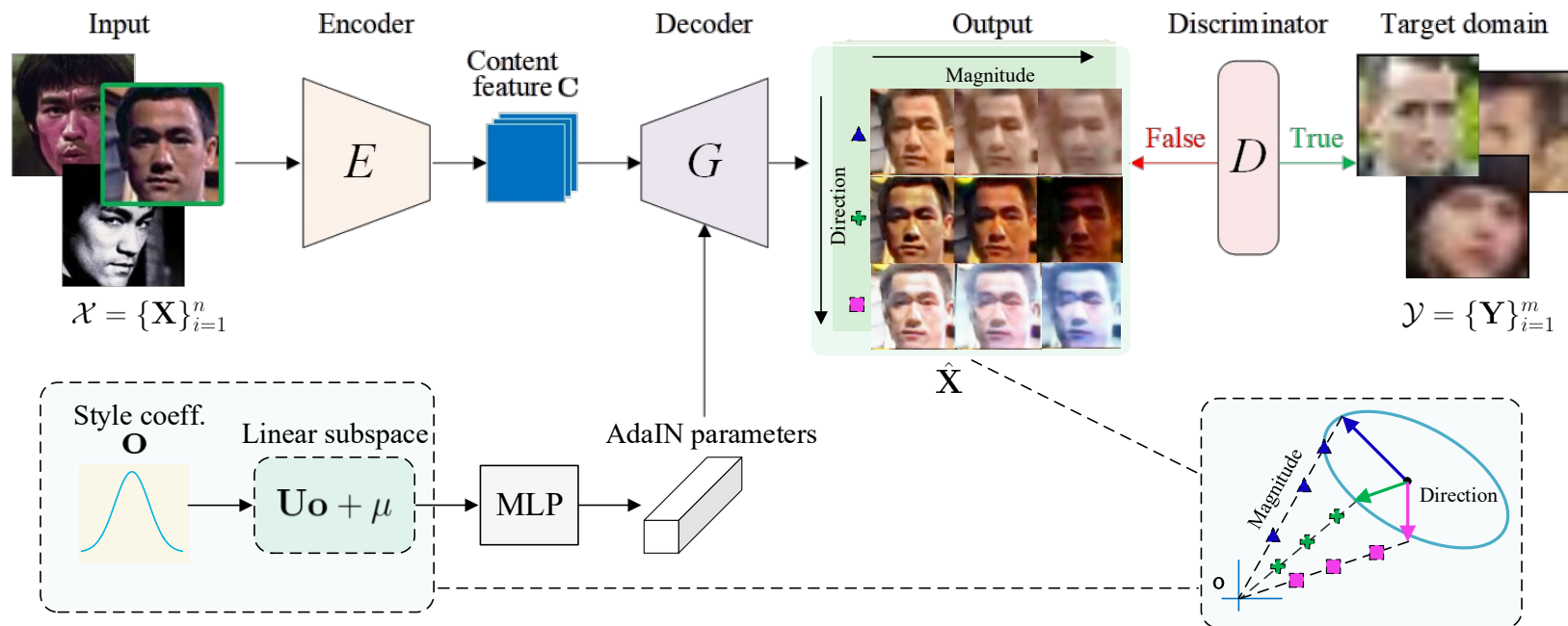- Bad illumination
- Turbulence effect

......

Synthesis Module

# Motivation

- Previous synthesis models: ***limited face properties***; offline and ***blind*** data augmentation.

- Facial properties should be generalizable to the challenging unconstrained testing scenarios.

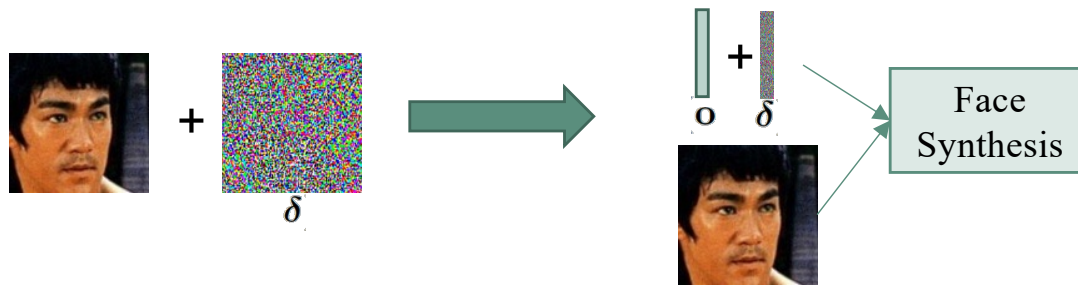- Feedback-based face synthesis is more beneficial to FR models.



- Low resolution
- Motion blurring
- Bad illumination
- Turbulence effect
……

Feedback

Synthesis Module

Recognition Module

# Controllable Face Synthesis

➢ The synthesis model can discover the **styles** in the target unconstrained data.

➢ The synthesis model is precisely-controllable in the style latent space, in both ***diversity*** and ***degree***.



Feng Liu, Minchul Kim, Anil Jain, and Xiaoming Liu. Controllable and Guided Face Synthesis for Unconstrained Face Recognition. ECCV 2022

# Guided Face Synthesis for Face Recognition

➢ The FR model feedback signal is incorporated into the face generation using *adversarial perturbation*.

➢ The manipulation of the low-dimensional style space renders this feedback *meaningful* and *efficient*.



- Style latent perturbations to maximize the classifier loss

$$\delta^* = \underset{||\delta|| < \epsilon}{\arg\max} \mathcal{L}_{cla}\left(\mathcal{F}(\mathbf{X}^*), l\right), \text{where } \mathbf{X}^* = G(E(\mathbf{X}), \text{MLP}(\mathbf{U}(\mathbf{o} + \boldsymbol{\delta}) + \boldsymbol{\mu}))$$
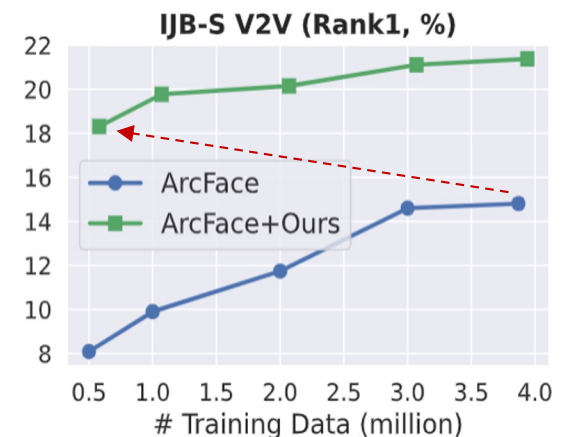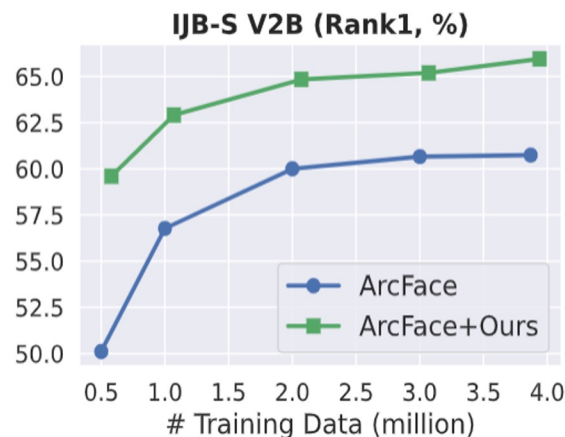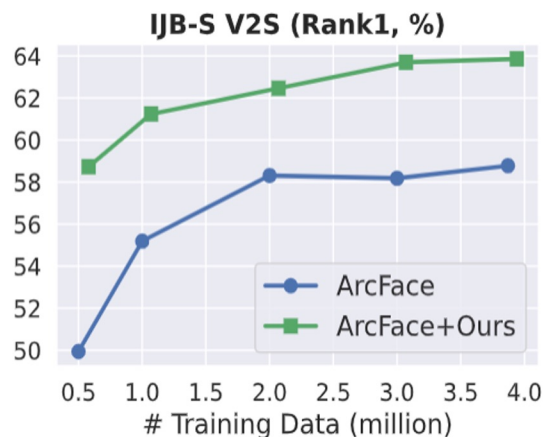
- Optimize the face embedding model

$$\min_{\theta} \mathcal{L}_{cla}([\mathbf{X}^*, \mathbf{X}], l)$$

Feng Liu, Minchul Kim, Anil Jain, and Xiaoming Liu. Controllable and Guided Face Synthesis for Unconstrained Face Recognition. ECCV 2022

# Face Recognition Results on IJB-S and TinyFace

➤ Our synthesis models could be plugged into any SoTA FR model and improve its performance.

| Method | Labeled Train Data | Backbone | IJB-S V2S | | | | IJB-S V2B | | | | IJB-S V2V | | | | TinyFace | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Rank1 | Rank5 | 1% | 10% | Rank1 | Rank5 | 1% | 10% | Rank1 | Rank5 | 1% | 10% | Rank1 | Rank5 |
| ArcFace[1] | MS1MV2-* | ResNet-50 | 58.78 | 66.40 | 40.99 | 50.45 | 60.66 | 67.43 | 43.12 | 51.38 | 14.81 | 26.72 | 2.51 | 5.72 | 62.21 | 66.85 |
| ArcFace+Ours* | MS1MV2-* | ResNet-50 | 61.69 | 68.33 | 43.99 | 53.34 | 62.20 | 69.50 | 44.38 | 53.49 | 18.14 | 31.34 | 2.09 | 4.51 | 62.39 | 67.36 |
| ArcFace+Ours | MS1MV2-* | ResNet-50 | **63.86** | **69.95** | **47.86** | **56.44** | **65.95** | **71.16** | **47.28** | **57.24** | **21.38** | **35.11** | **2.96** | **7.41** | **63.01** | **68.21** |
| AdaFace[2] | WebFace12M | IResNet-100 | 71.35 | 76.24 | 59.40 | **66.34** | 71.93 | 76.56 | 59.37 | **66.68** | 36.71 | 50.03 | 4.62 | 11.84 | 72.29 | 74.97 |
| AdaFace+Ours | WebFace12M | IResNet-100 | **72.54** | **77.59** | **60.94** | 66.02 | **72.65** | **78.18** | **60.26** | 65.88 | **39.14** | **50.91** | **5.05** | **13.17** | **73.87** | **76.77** |

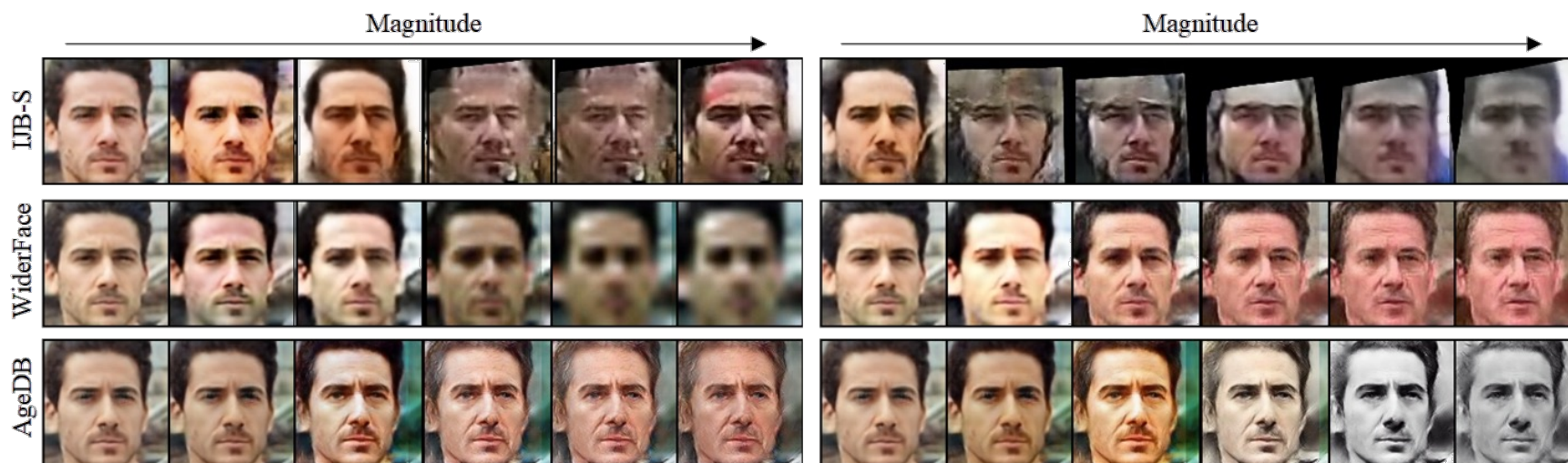➤ Our synthesis models can boost FR performance even with less labelled training samples.



[1] Arcface: Additive angular margin loss for deep face recognition. CVPR 2019.
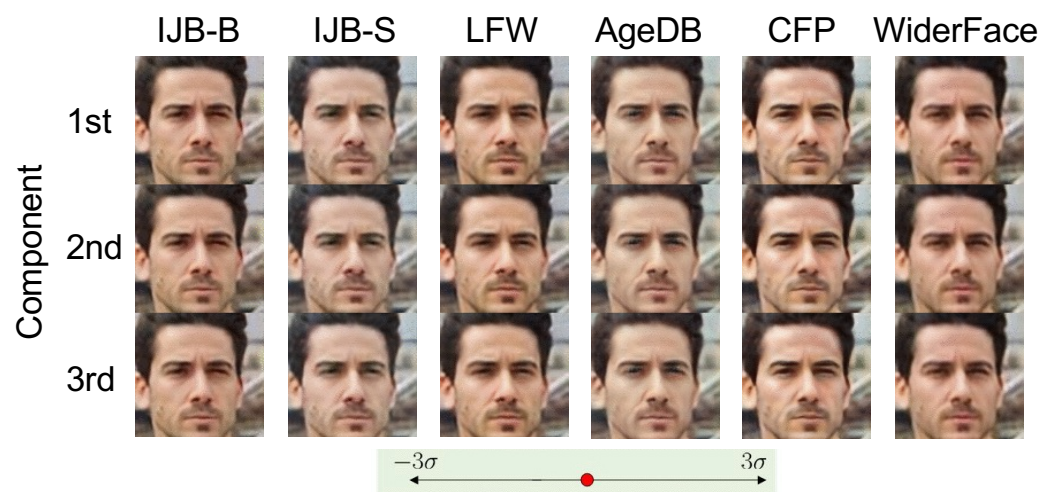[2] AdaFace: Quality Adaptive Margin for Face Recognition. CVPR 2022.

# Visualizations of the Face Synthesis Model

➢ Interpretable magnitude of the style coefficient.



➢ Learned the orthonormal basis of the subspace.

➢ The distribution similarity between datasets A and B

$$\mathcal{S}(A, B) = \frac{1}{q} \left( \sum_i^q S_C(\mathbf{u}_A^i + \boldsymbol{\mu}_A, \mathbf{u}_B^i + \boldsymbol{\mu}_B) \right)$$
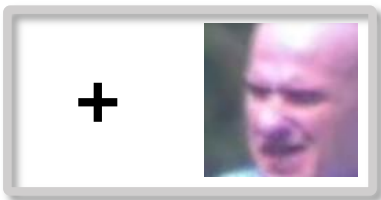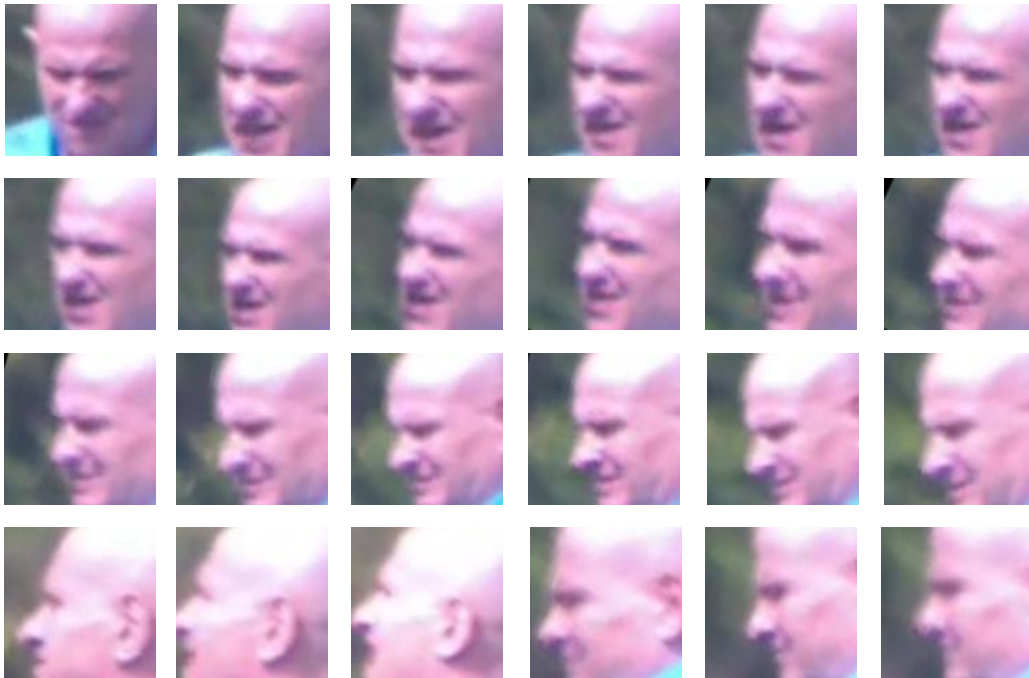
# Cluster and Aggregate:
# Face Recognition with Large Probe Set

Minchul Kim, Feng Liu, Anil K. Jain, Xiaoming Liu
NeurIPS 2022

# Problem Definition

## Traits of Face Recognition with Videos



**1** Varied Identifiability

Some images are more identifiable than others

**2** 10 ~ 1,000,000

Varied Number of Images

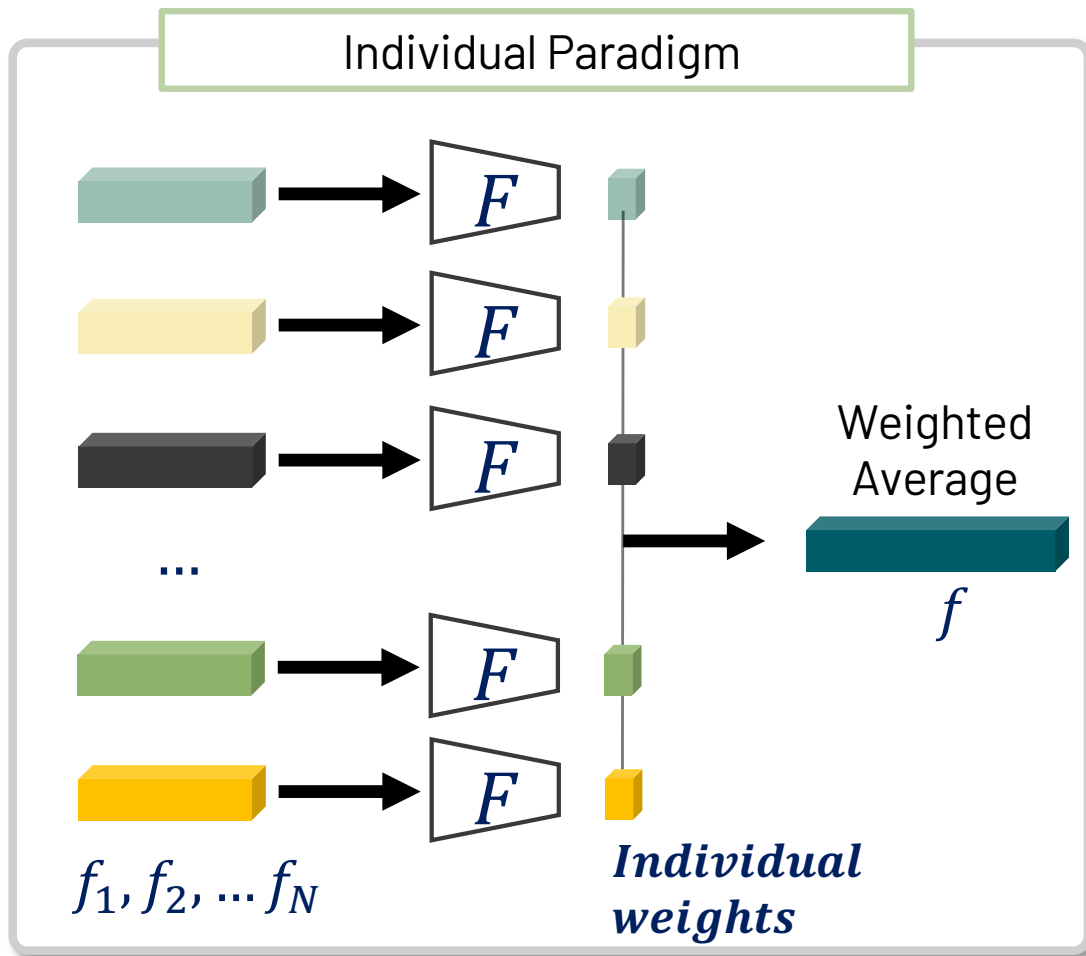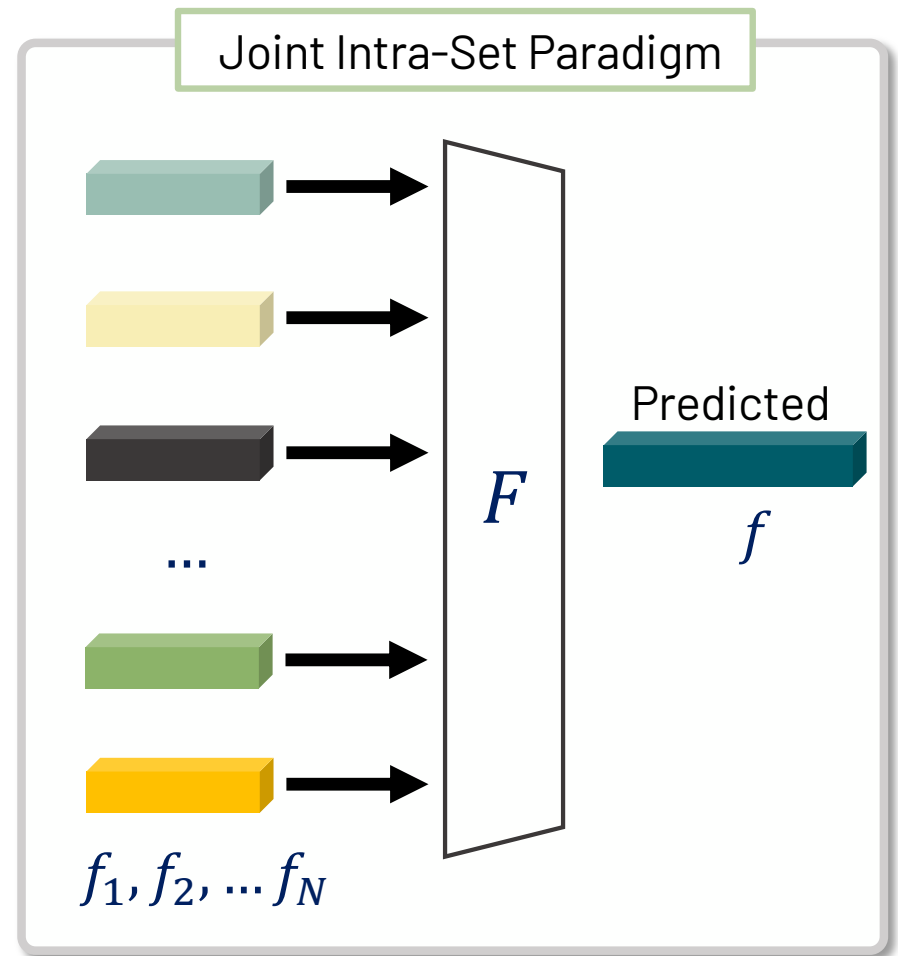Number of images are not fixed.

**3** Sequential Inputs

Videos come in sequentially. We use what we have up-to the current timeframe.

Subject in the slide consented to publication
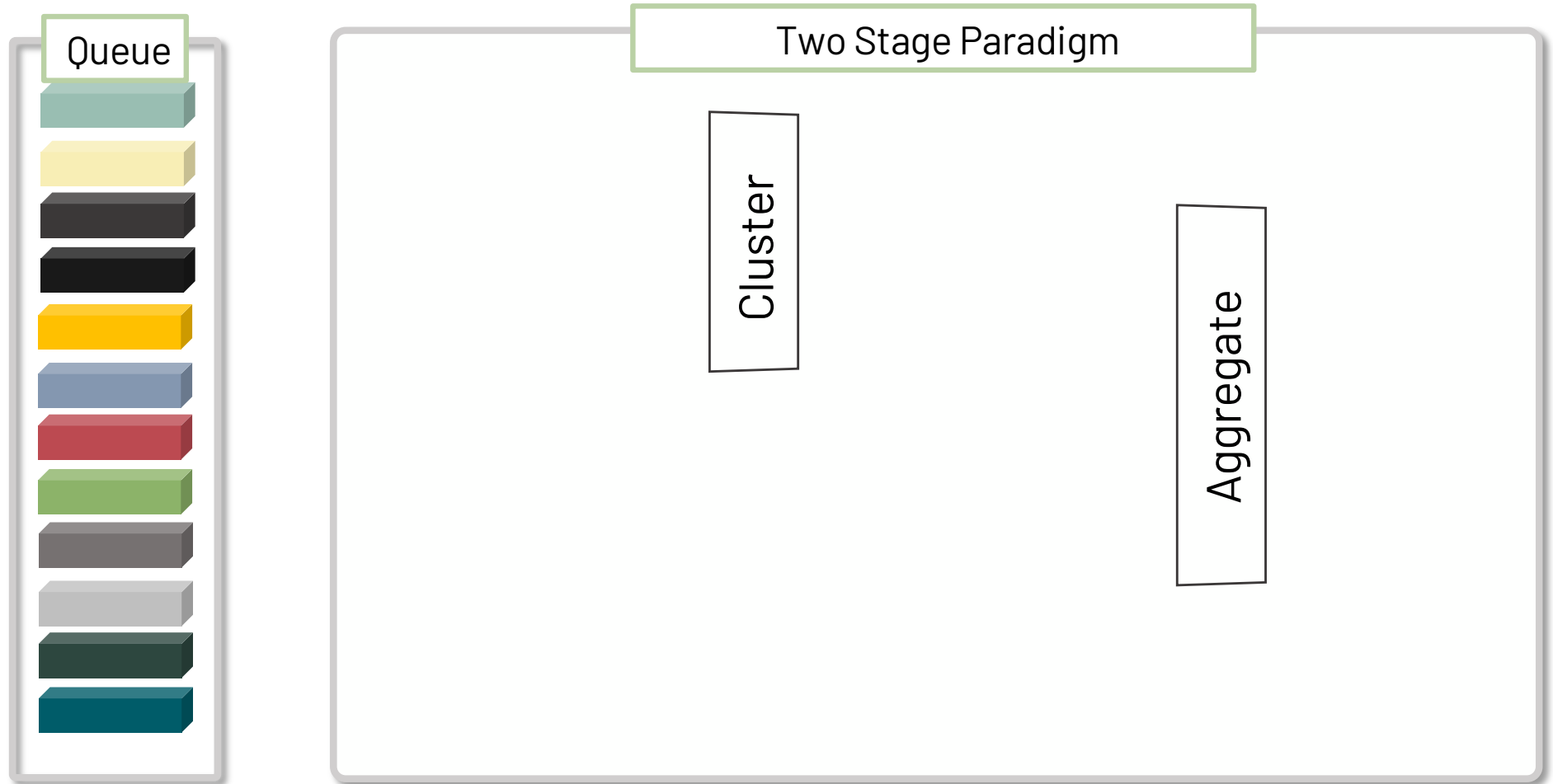
# Problem of Previous Methods



Individual Paradigm

$f_1, f_2, \dots f_N$

$F$

Individual weights

Weighted Average

$f$

(No Intra−Set Relationship)

Joint Intra-Set Paradigm

$f_1, f_2, \dots f_N$

$F$

Predicted

$f$

(Cannot handle large N)

## Large N / Sequential Scenario

Queue

Two Stage Paradigm

Cluster

Aggregate

# Motivation

## Large N / Sequential Scenario

## Large N / Sequential Scenario



Queue

Two Stage Paradigm

$T_2$

$f_5, f_6, f_7, f_8$

Cluster

$f'_1, f'_2, f'_3$

Summarized into
3 features

Aggregate

$f$

## Large N / Sequential Scenario

# Method

## Architecture



**Cluster and Aggregate (CAFace)**

$N'$ images

Style Input Maker (SIM)

$E$ — Decouples style and identity

Style $S$ — $\mathbb{R}^{N' \times d}$

Identity $F$ — $\mathbb{R}^{N' \times C}$

**Cluster Network (CN)**

$C \in \mathbb{R}^{M \times d}$

Global Centers

Key: $S$     Query: $C$

Transformer + SoftMax

**Assignment map $A$**   $M \times N'$

Value: $(S, F)$

Assigns inputs to globally shared centers

$S'_{t-1}, F'_{t-1}, A_{t-1}$ — previous batch

$S' = AS$ — $\mathbb{R}^{M \times d}$

$F' = AF$ — $\mathbb{R}^{M \times C}$

$S'_t, F'_t, A_t$ — next batch

**Aggregation Network (AGN)**

$P \in \mathbb{R}^{M \times C}$

MLP MIxer

weight

weighted average

$f$: output — $\mathbb{R}^{1 \times C}$

(Optionally) update intermediates with previous batch intermediates and fuse.

---
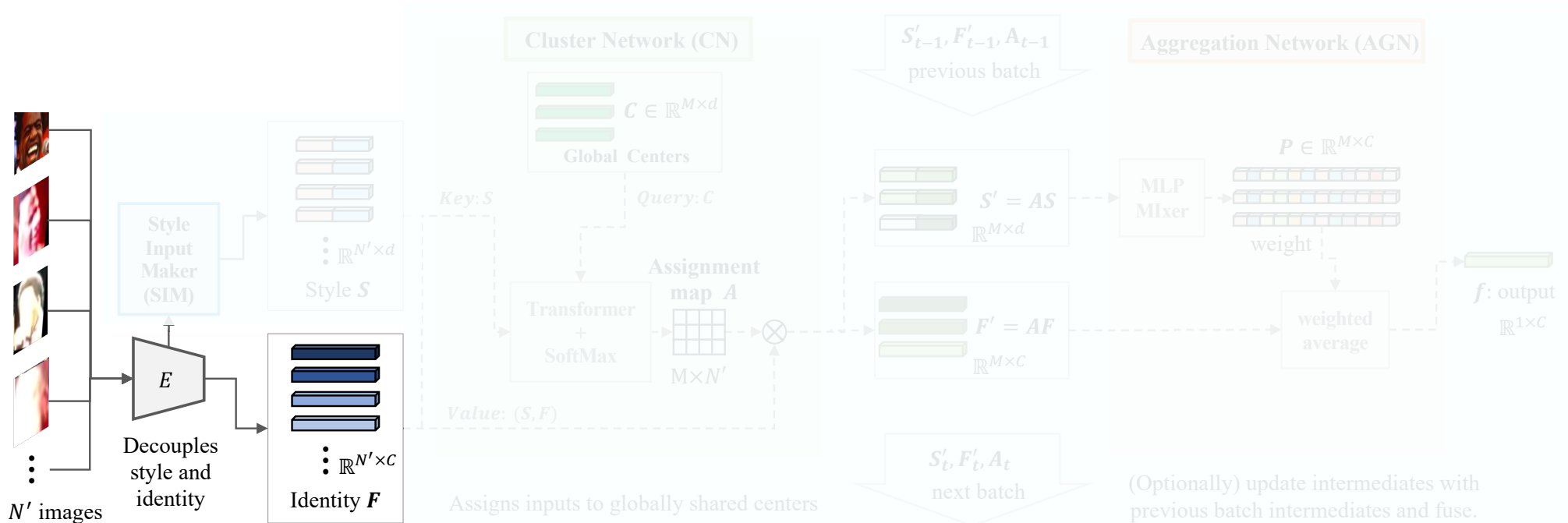
Overall Architectures
3 components (SIM, CN, AGN)

## Architecture



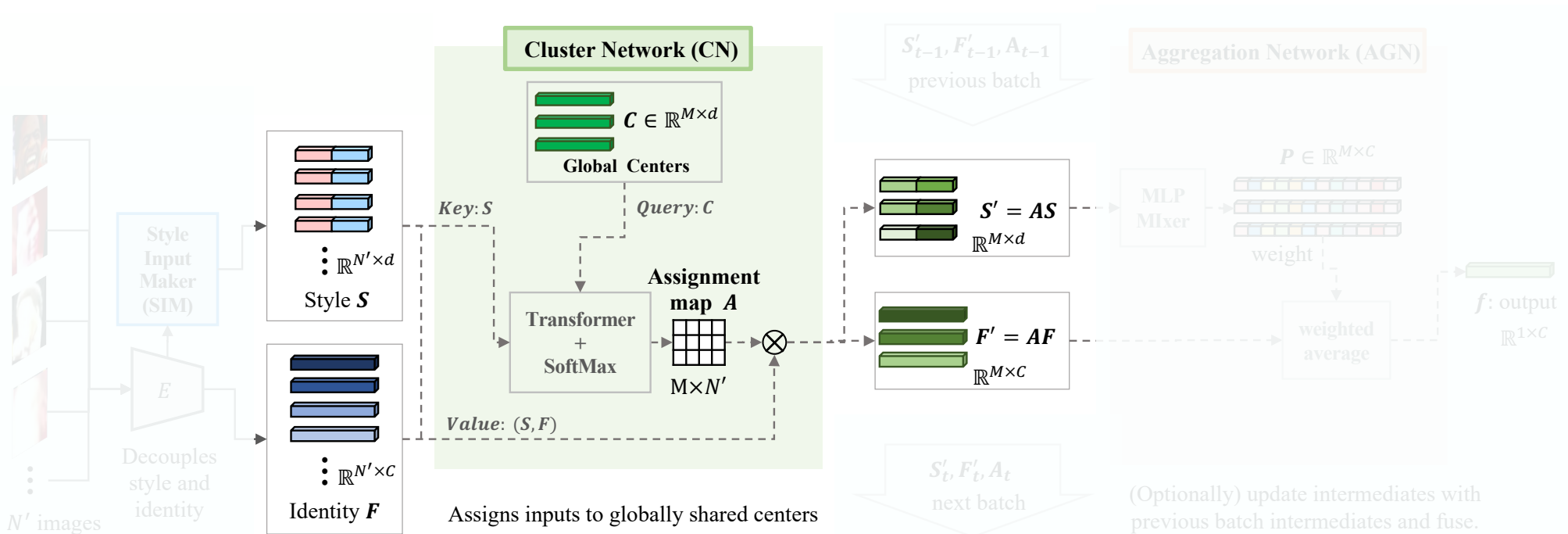Input images fed into the fixed feature extractor.
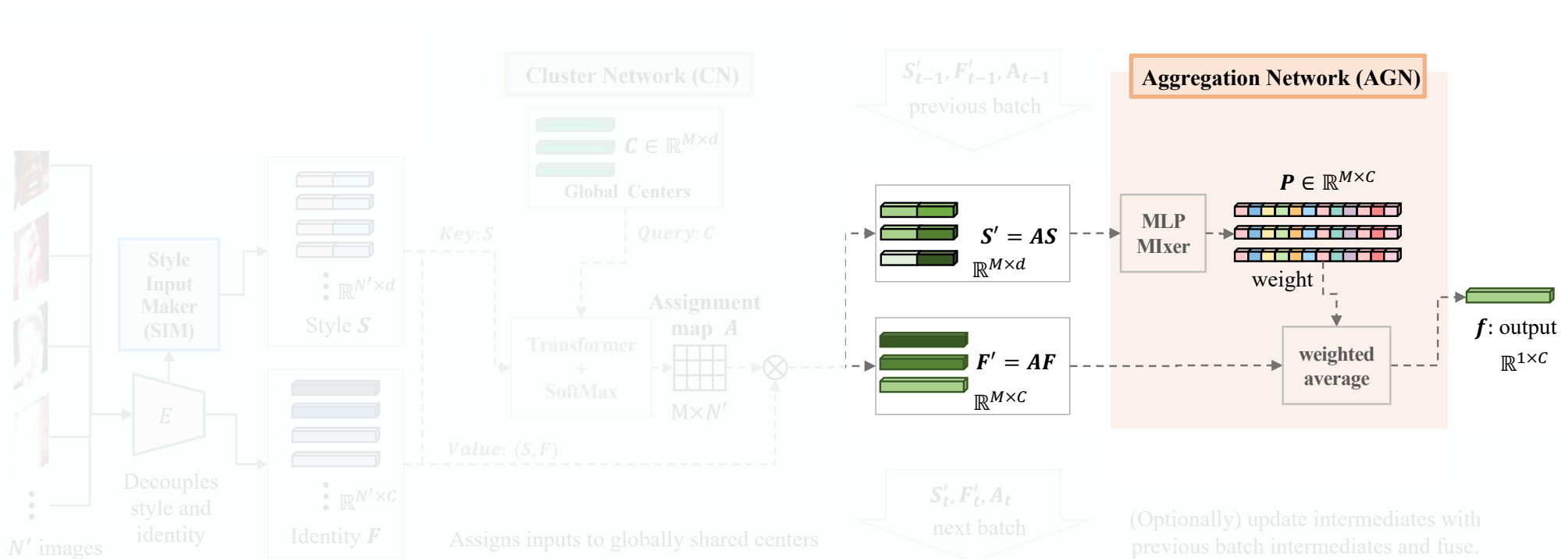
# Method

## Architecture

# Method

## Architecture

CN uses learned centers $\{c_j\}^M$ and $\{s_i\}^N$ to create assignment map $A$.
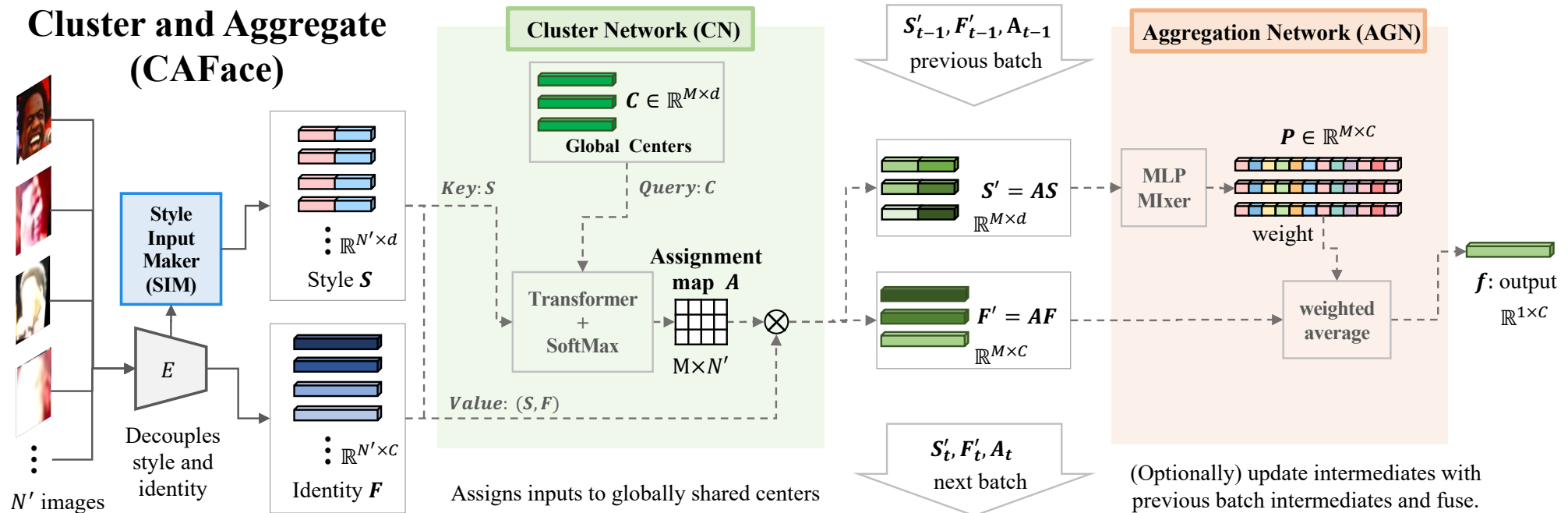$A$ is used to map $\{f_i\}^N \to \{f_j\}^M$ and $\{s_i\}^N \to \{s_j\}^M$

# Method

## Architecture



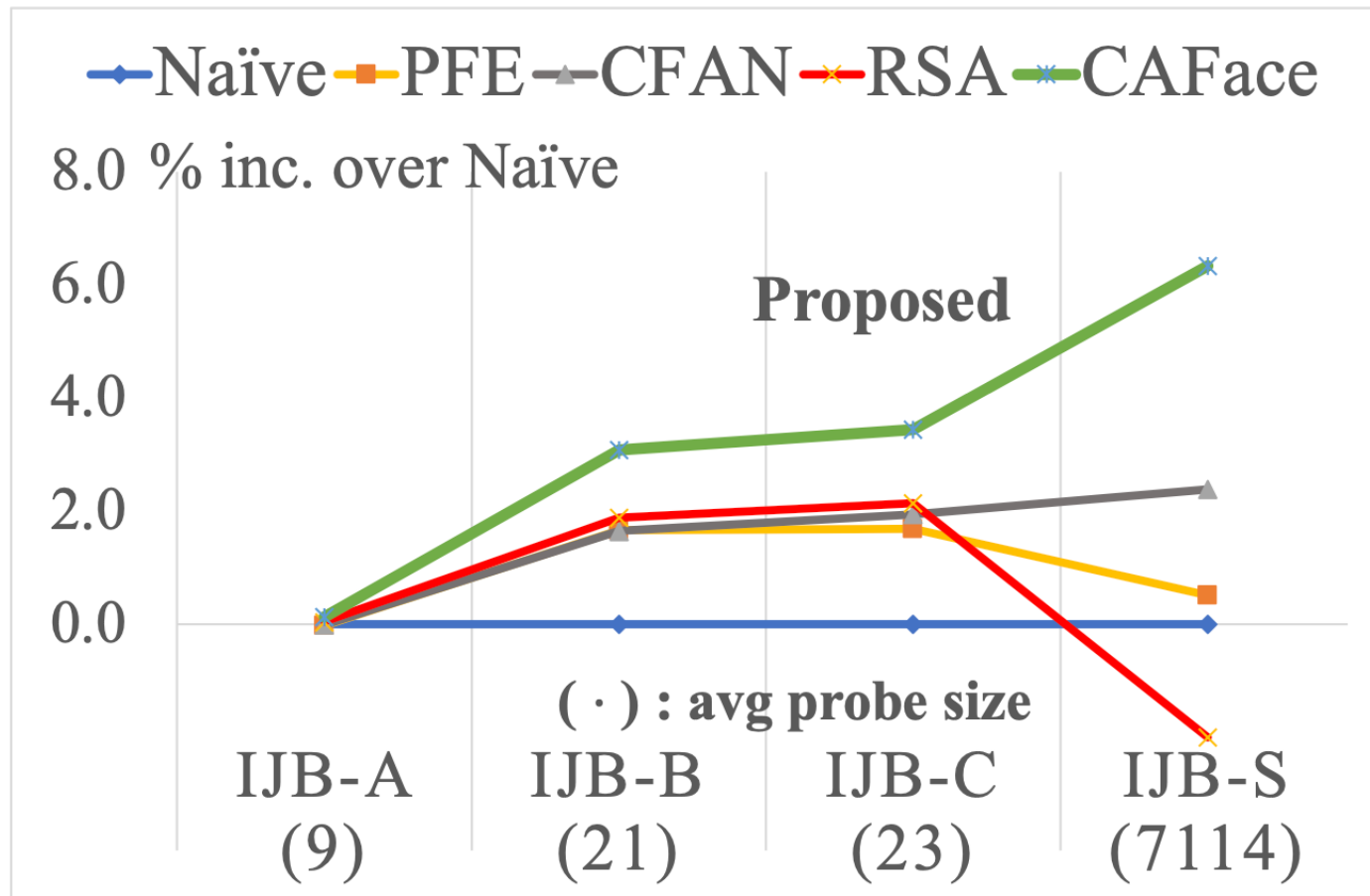AGN maps $\{f_j\}^M, \{s_j\}^M \rightarrow f$ with intra-set relationship.

## Architecture



**Cluster and Aggregate (CAFace)**

$N'$ images

Decouples style and identity

Style Input Maker (SIM)

$E$

Style $S$ $\mathbb{R}^{N' \times d}$

Identity $F$ $\mathbb{R}^{N' \times C}$

**Cluster Network (CN)**

$C \in \mathbb{R}^{M \times d}$

Global Centers

Key: $S$    Query: $C$

Transformer + SoftMax

**Assignment map $A$**

$M \times N'$

Value: $(S, F)$

Assigns inputs to globally shared centers

$S'_{t-1}, F'_{t-1}, A_{t-1}$

previous batch

$S' = AS$ $\mathbb{R}^{M \times d}$

$F' = AF$ $\mathbb{R}^{M \times C}$

$S'_t, F'_t, A_t$

next batch

**Aggregation Network (AGN)**

$P \in \mathbb{R}^{M \times C}$

MLP MIxer

weight

weighted average

$f$: output $\mathbb{R}^{1 \times C}$

(Optionally) update intermediates with previous batch intermediates and fuse.

Intermediate features $\{f_j\}^M$ and $\{s_j\}^M$ are updated in sequential setting.

**Performance Gain over simple average using feature fusion methods.**



**Naïve**: Simple Average
**PFE, CFAN**: single image weight estimation
**RSA**: Attention Mechanism

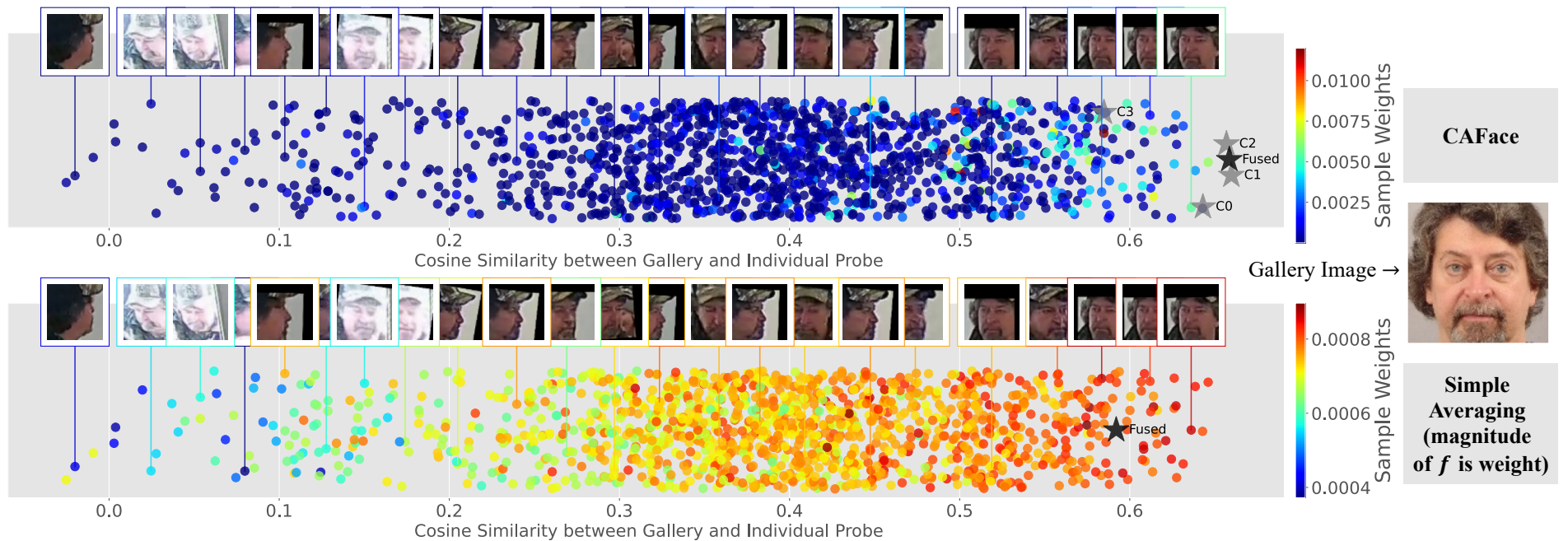**Largest Probe size, Largest Perf. gain**

**Table 3:** A performance comparison of recent methods on the IJB-S [24] dataset.

| Method | Surveillance-to-Single | | | Surveillance-to-Booking | | | Surveillance-to-Surveillance | | |
|---|---|---|---|---|---|---|---|---|---|
| | Rank-1 | Rank-5 | 1% | Rank-1 | Rank-5 | 1% | Rank-1 | Rank-5 | 1% |
| Naive Average | 69.26 | 74.31 | 57.06 | 70.32 | 75.16 | 56.89 | 32.13 | 46.67 | 5.32 |
| PFE [46] | 69.50 | 74.39 | 57.51 | 70.53 | 75.29 | 57.98 | 32.27 | 46.70 | 5.41 |
| CFAN [15] | 70.00 | 74.58 | 57.93 | 70.90 | 75.58 | 58.09 | 31.66 | 45.59 | 5.79 |
| RSA [31] | 63.04 | 67.33 | 51.62 | 63.54 | 68.23 | 51.89 | 16.82 | 31.80 | 0.75 |
| CAFace | **71.61** | **76.43** | **62.21** | **72.72** | **77.41** | **62.68** | **36.51** | **49.59** | **8.78** |
| CAFace (Random Order) | 71.65 ±0.05 | 76.37 ±0.04 | 62.27 ±0.11 | 72.77 ±0.04 | 77.37 ±0.03 | 62.70 ±0.06 | 36.43 ±0.08 | 49.40 ±0.05 | 8.89 ±0.03 |

Changing the order of probe sequence does not affect the performance.

# Visualization of Assignment Map

Assignment Map $A$ Visualization

Each cluster is formed by weighted averaging each row.



| Cluster 1 | Mean $P_1$ : 0.653 | 0.45 | .24 | 0.21 | 0.20 | 0.20 | 0.10 | 0.04 | 0.04 | 0.03 | 0.03 | 0.02 | 0.01 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| Cluster 2 | Mean $P_2$ : 0.258 | 0.30 | .41 | 0.25 | 0.25 | 0.29 | 0.09 | 0.24 | 0.08 | 0.08 | 0.06 | 0.16 | 0.01 | 0.07 | 0.09 | 0.04 | 0.00 | 0.01 | 0.02 | 0.02 | 0.00 | 0.00 | 0.00 | 0.00 |
| Cluster 3 | Mean $P_3$ : 0.089 | 0.22 | .18 | 0.33 | 0.35 | 0.30 | 0.53 | 0.07 | 0.42 | 0.21 | 0.38 | 0.14 | 0.19 | 0.07 | 0.01 | 0.02 | 0.09 | 0.03 | 0.01 | 0.00 | 0.01 | 0.00 | 0.00 | 0.00 |
| Cluster 4 | Mean $P_4$ : 0.000 | 0.04 | .18 | 0.21 | 0.20 | 0.22 | 0.28 | 0.65 | 0.46 | 0.68 | 0.53 | 0.68 | 0.79 | 0.86 | 0.90 | 0.94 | 0.90 | 0.96 | 0.97 | 0.97 | 0.99 | 1.00 | 1.00 | 1.00 |
| $P \in \mathbb{R}^{4 \times 512}$ | Mean $\mathbf{P}_j$ : $\mathbb{R}^{512} \rightarrow \mathbb{R}^1$ | | | | | | | | | | | | | | | | | | | | | | | |

Importance of each cluster during aggregation.

Each Column sums up to 1. They are soft assigned to cluster centers.

# Weight Visualizations



IJBS Probes' similarity to Gallery Visualization

Point colors indicate the weight during fusion.

# CAFace Demo

**Video Feed - Frame Count: 1**

**Detected Face Frame**

**Top 5 CAFace Match**

| |
|---|
| 0.13 |
| 0.12 |
| 0.10 |
| 0.10 |
| 0.10 |

**Top 5 Naive Match**

| |
|---|
| 0.13 |
| 0.12 |
| 0.10 |
| 0.10 |
| 0.10 |

**Tok 5 Most influential Images in CAFace**

**Cosine Similarity To Gallery**

- CAFace
- Naive

**Gallery Image**

-0.3  -0.2  -0.1  0.0  0.1  0.2  0.3  0.4  0.5  0.6  0.7  0.8  0.9  1.0

# CAFace Demo

**Video Feed - Frame Count: 1**



**Detected Face Frame**

**Top 5 CAFace Match**

0.13
0.12
0.12
0.11
0.11

**Top 5 Naive Match**

0.13
0.12
0.12
0.11
0.11

**Tok 5 Most influential Images in CAFace**

**Cosine Similarity To Gallery**

**Gallery Image**

CAFace
Naive

-0.3  -0.2  -0.1  0.0  0.1  0.2  0.3  0.4  0.5  0.6  0.7  0.8  0.9  1.0

# People Matching: Learning Clothing Invariant 3D Shape Representation

Feng Liu, Minchul Kim, ZiAng Gu, Anil Jain, and Xiaoming Liu
Under review

People matching. Two main characteristics: diverse human activities and clothing changes



Gait recognition

Person re-identification

People matching

Disentangle identity and non-identity features in 3D body shape space

## Disentangle identity and non-identity features in 3D body shape space

Disentangle identity and non–identity features in 3D body shape space

# Diverse People Matching Dataset (DPMD)

87,821 images of 536 subjects

Examples of diverse poses

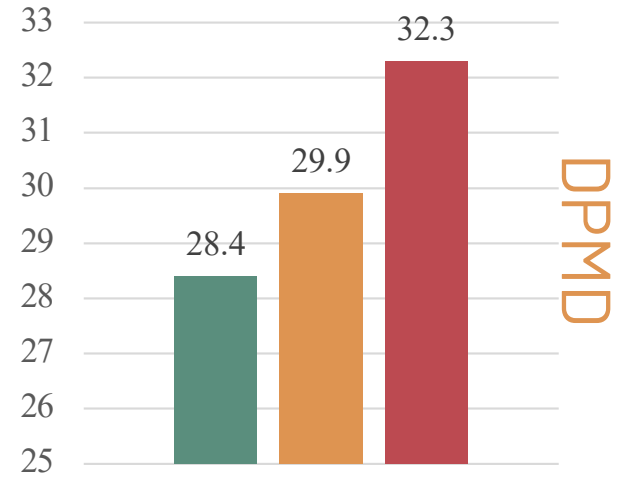Examples of diverse clothes

# People Matching Results
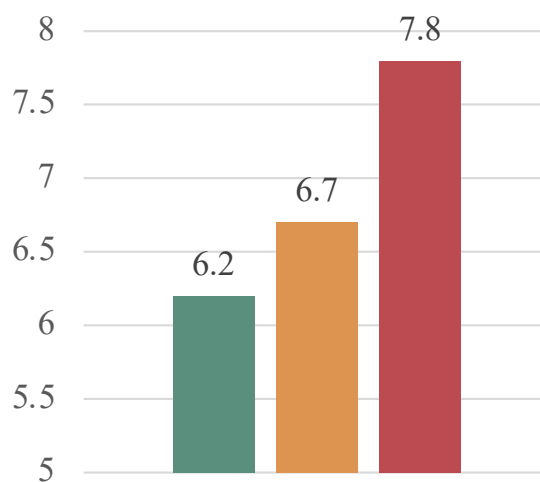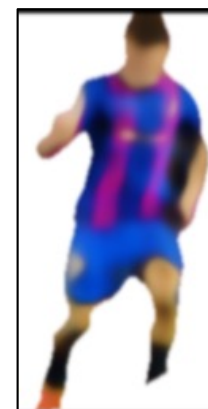
# 3D Reconstruction Results



Naked body     Clothed body     Rec. image         Naked body     Clothed body     Rec. image

# Conclusions

➤ There are many research questions for low-quality recognition

➤ Even for conventional FR problems, there are research opportunities such as explainability, new architecture, etc.

➤ Body biometrics is just at the beginning and there is a great potential for further development.

# Questions?