國立高雄科技大學
National Kaohsiung University of Science and Technology

# FaceNet: A Unified Embedding for Face Recognition and Clustering
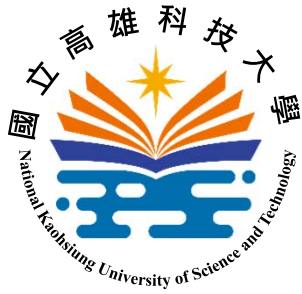
Florian Schroff, Dmitry Kalenichenko, and James Philbin

Speaker: Shih-Shinh Huang
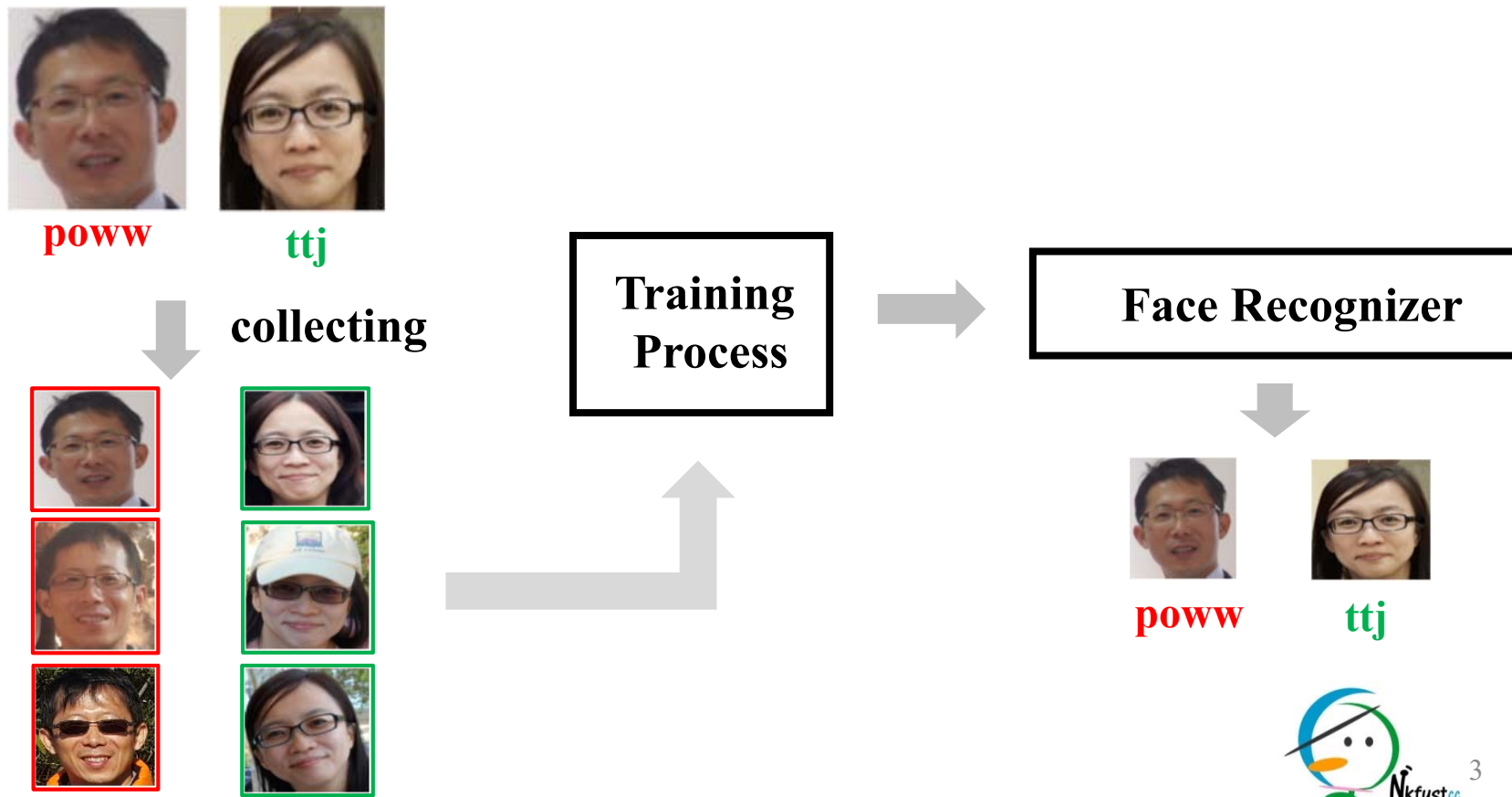
March 5, 2020

# Outline

- Introduction
  - Motivation
  - Idea of FaceNet
  - Benefit from FaceNet

- Euclidean Embedding
  - FaceNet Architecture
  - Triplet Loss

- Triplet Selection
  - Motivation
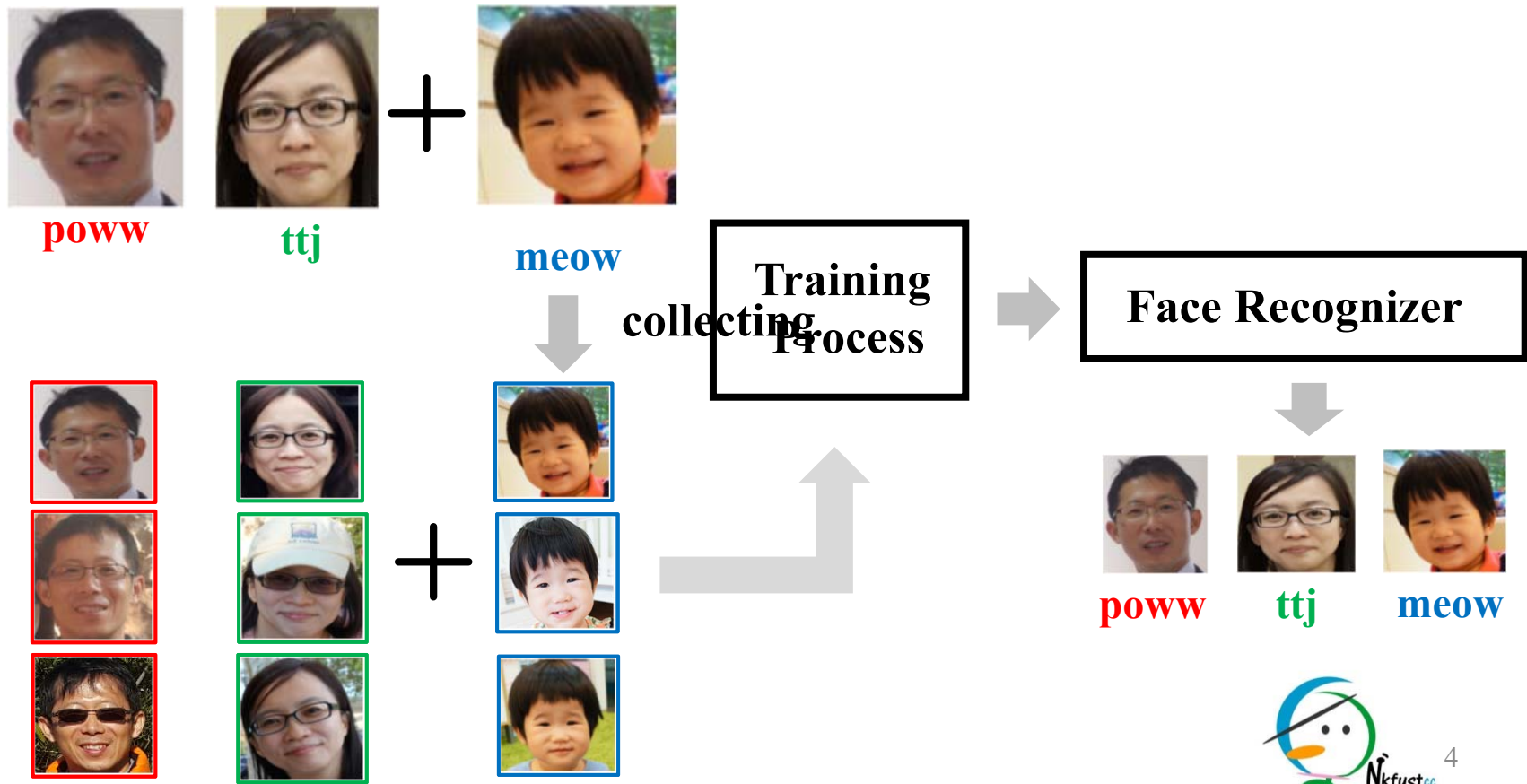  - Intuitive Selection
  - Online Generation

# Introduction

- Motivation: Previous Face Recognition

# Introduction

- Motivation: Previous Face Recognition
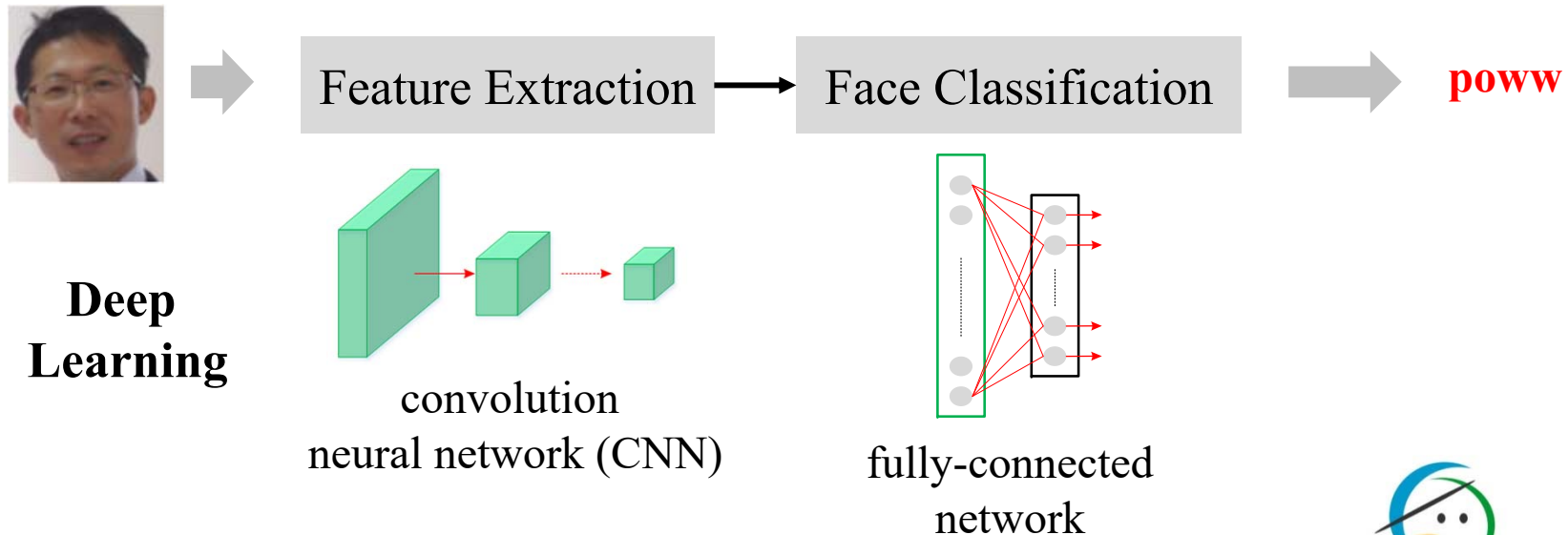
# Introduction

- Motivation: Problems
  - Retraining Problem: the cost of re-training is high and is sometimes not affordable.
  - Data Collection Problem: collecting sufficient faces to training a good classifier is impractical.

# Introduction

- Idea of FaceNet
  - A face recognizer generally consists of two stages
    - Feature Extraction: describe a face in an effective way.
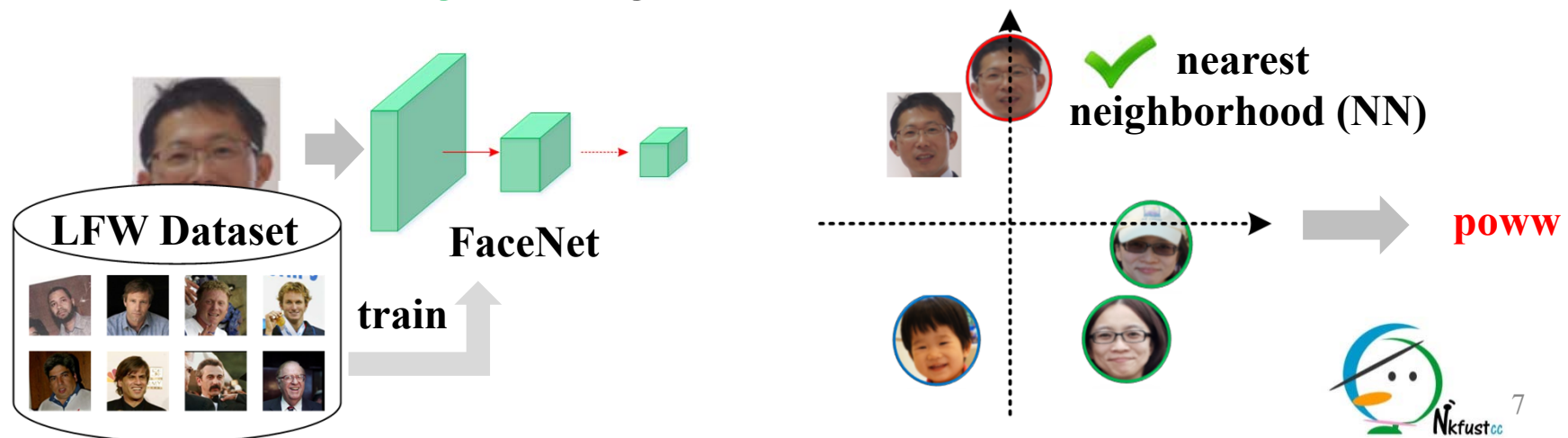    - Classification: assign an identifier to a face



**Deep Learning**

Feature Extraction → Face Classification → **poww**

convolution neural network (CNN)

fully-connected network

- Idea of FaceNet
  - de-couple two face recognition stages

Feature Extraction ❌ Face Classification

- FaceNet is CNN just for feature using simple but not training classification.
  - Face Net is a CNN achieved by using simple but not training algorithms.



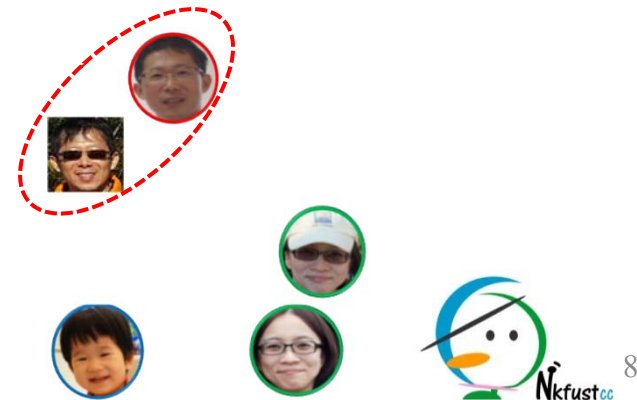LFW Dataset

FaceNet

train

nearest neighborhood (NN)

poww

# Introduction

- Benefit from FaceNet
  - The trained FaceNet is applicable to any face applications without re-training.
  - The use of FaceNet makes the number of faces necessary for face-related tasks small

verification → thresholding

classification → nearest neighborhood



$d_{th}$
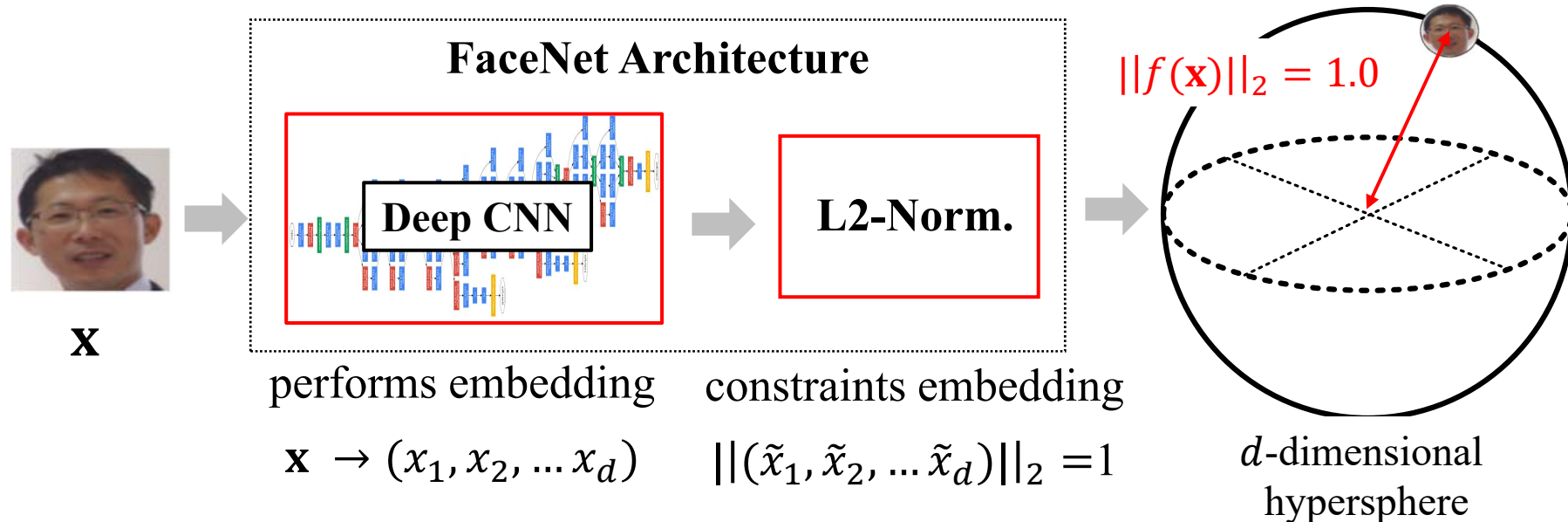
# Euclidean Embedding

- FaceNet Architecture $f(.)$
  - maps a face clip to Euclidean space
  - makes the similarity measure possible.

$$\sqrt{(0.7-(-1.0))^2 + (0.4-(-0.9))^2}$$

Euclidean Embedding (FaceNet)

$x_1$: **hair length**

$x_2$: **skin color**

short   long   $x_1$

$\begin{pmatrix} +0.7 \\ +0.4 \end{pmatrix}$

$\begin{pmatrix} -1.0 \\ -0.9 \end{pmatrix}$

dark

# Euclidean Embedding

- FaceNet Architecture $f(.)$
  - Formal Description: maps a face $\mathbf{x}$ to a $d$-dimensional unit vector $f(\mathbf{x}) = (\tilde{x}_1, \tilde{x}_2, ... \tilde{x}_d)$

**FaceNet Architecture**

**Deep CNN**

**L2-Norm.**

$\mathbf{x}$

$||f(\mathbf{x})||_2 = 1.0$

performs embedding

$\mathbf{x} \rightarrow (x_1, x_2, ... x_d)$

constraints embedding

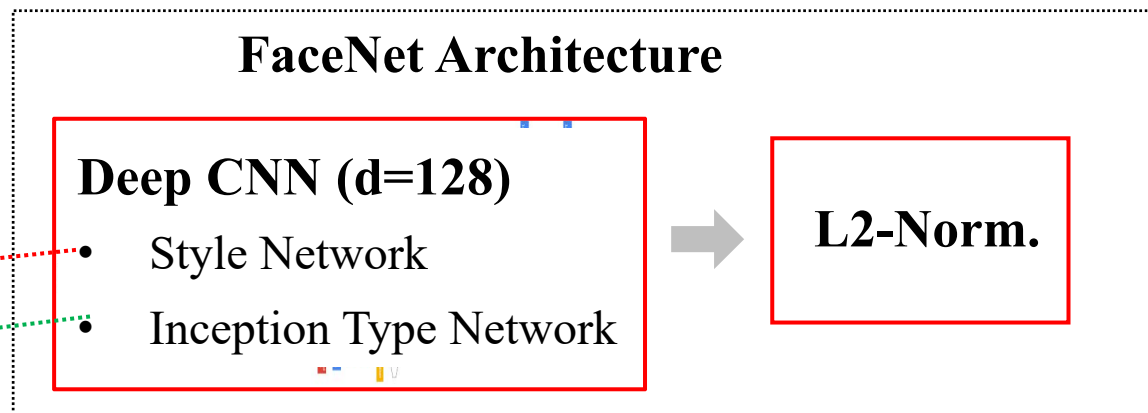$||(\tilde{x}_1, \tilde{x}_2, ... \tilde{x}_d)||_2 = 1$

$d$-dimensional hypersphere

# Euclidean Embedding
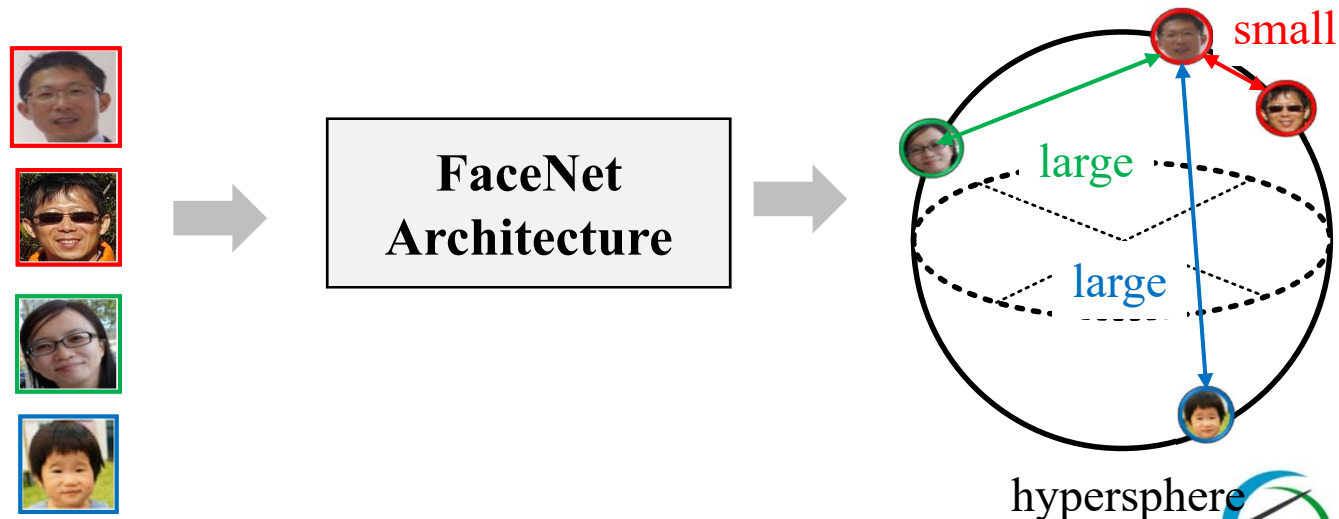
- FaceNet Architecture $f(.)$
  - Formal Description: maps a face $\mathbf{x}$ to a $d$-dimensional unit vector $f(\mathbf{x}) = (\tilde{x}_1, \tilde{x}_2, \ldots \tilde{x}_d)$

**FaceNet Architecture**

**Deep CNN (d=128)**
- Style Network
- Inception Type Network

→ **L2-Norm.**

- M. D. Zeiler and R. Fergus, "Visualizing and Understanding Convolutional Networks," CoPR, 2013
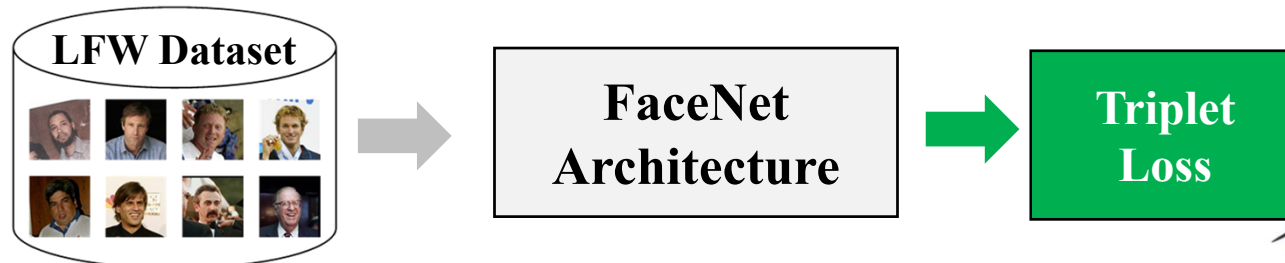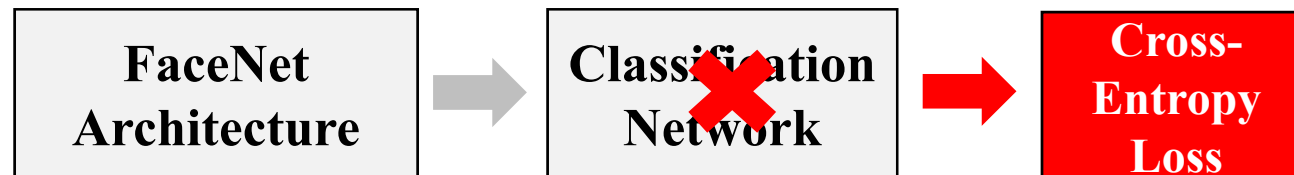- C. Szegedy, et. al. "Going Deeper with Convolutions," CoPR, 2014

# Euclidean Embedding

- FaceNet Architecture $f(.)$
  - Objective: provide discriminative embedding
    - faces of the same person have small distances
    - faces of distinct persons have large distances

# Euclidean Embedding

- Triplet Loss
  - It is a loss defined in embedding feature space.
    - be independent of faces to be recognized
    - be trained by using off-the-shelf face datasets

```
[FaceNet Architecture] ➡ [Classification ✖ Network] ➡ [Cross-Entropy Loss]
```

```
[LFW Dataset] ➡ [FaceNet Architecture] ➡ [Triplet Loss]
```

# Euclidean Embedding

- Triplet Loss
  - models embedding discriminability on a triplet $(\mathbf{x}^a, \mathbf{x}^p, \mathbf{x}^n)$

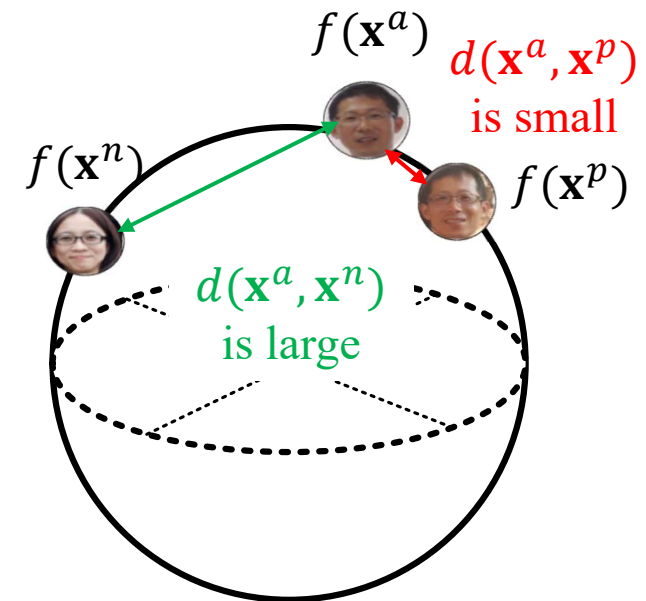$\mathbf{x}^a$
anchor image

$\mathbf{x}^p$
positive image

$\mathbf{x}^n$
negative image

**FaceNet Architecture**

$d(\mathbf{x}^a, \mathbf{x}^p) = \|f(\mathbf{x}^a) - f(\mathbf{x}^p)\|_2^2$

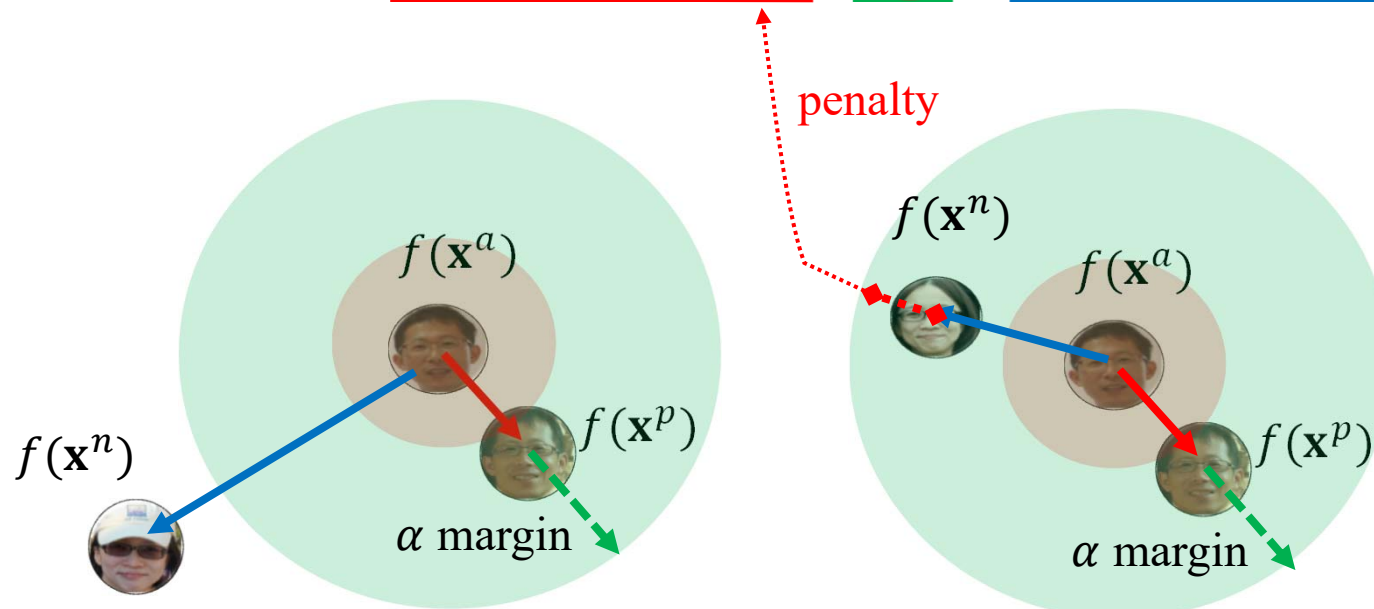$d(\mathbf{x}^a, \mathbf{x}^n) = \|f(\mathbf{x}^a) - f(\mathbf{x}^n)\|_2^2$

$f(\mathbf{x}^a)$

$d(\mathbf{x}^a, \mathbf{x}^p)$ is small

$f(\mathbf{x}^n)$

$f(\mathbf{x}^p)$

$d(\mathbf{x}^a, \mathbf{x}^n)$ is large

# Euclidean Embedding

- Triplet Loss

$$Loss(\mathbf{x}^a, \mathbf{x}^p, \mathbf{x}^n)$$

$$= \max\{d(f(\mathbf{x}^a), f(\mathbf{x}^p)) + \alpha - d(f(\mathbf{x}^a), f(\mathbf{x}^n)), \quad 0\}$$

penalty



$f(\mathbf{x}^a)$

$f(\mathbf{x}^p)$

$f(\mathbf{x}^n)$

$\alpha$ margin
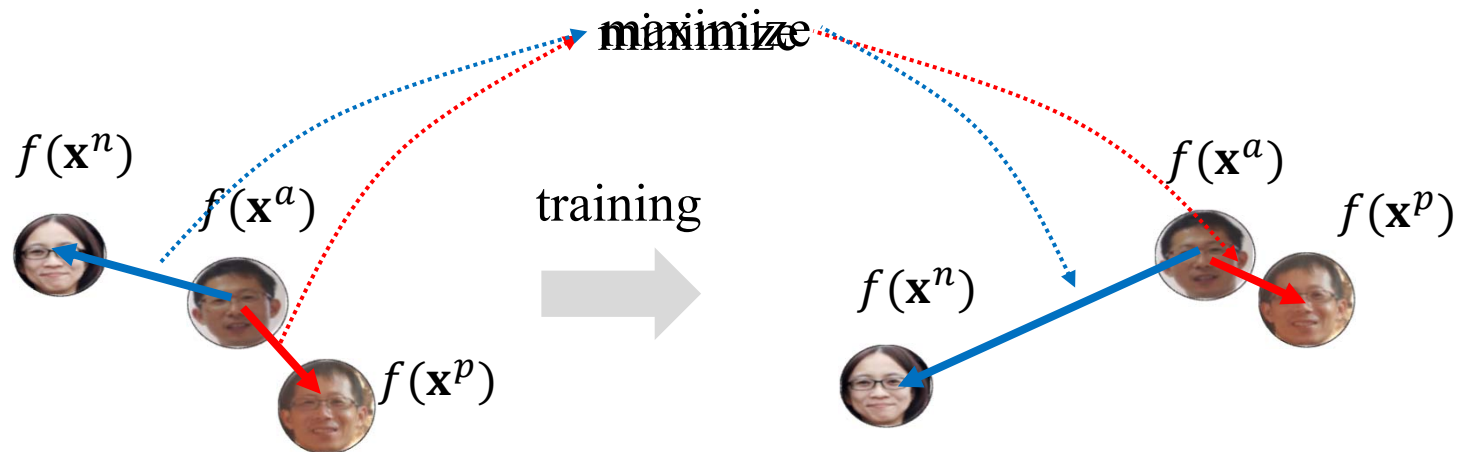
$f(\mathbf{x}^n)$

$f(\mathbf{x}^a)$

$f(\mathbf{x}^p)$

$\alpha$ margin

# Euclidean Embedding

- Triplet Loss

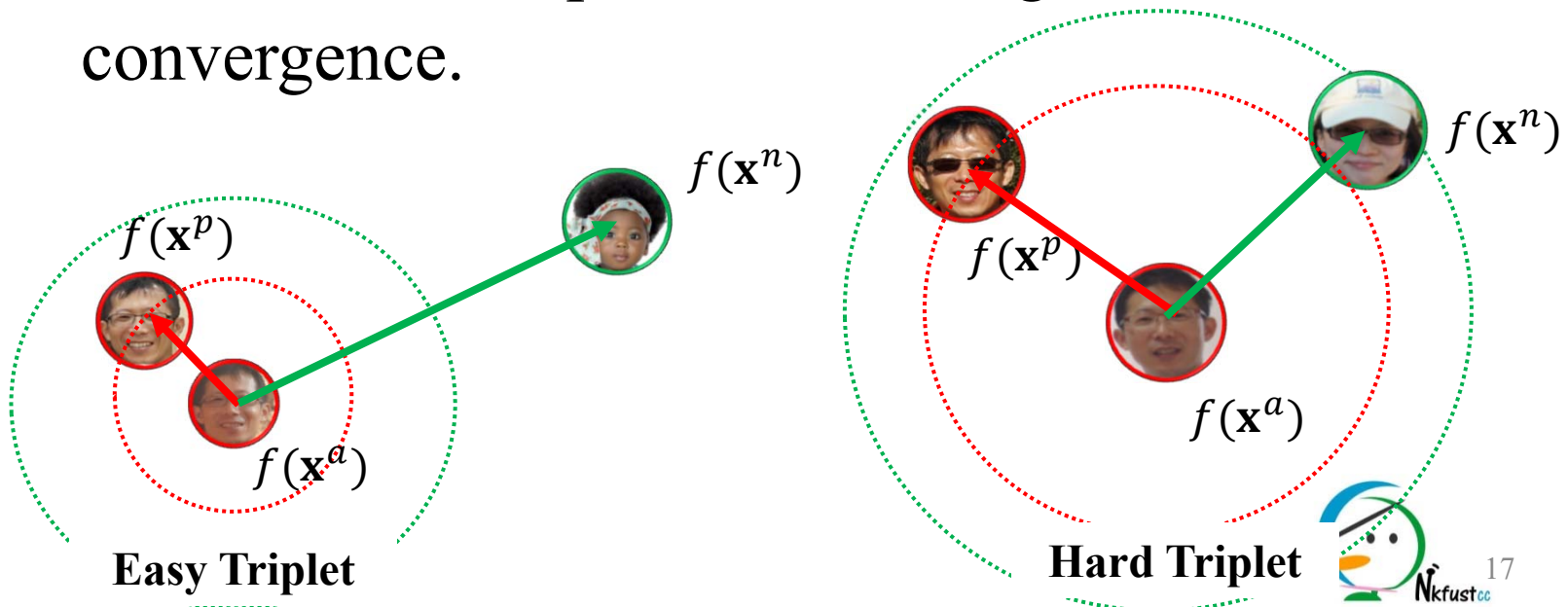$$L(T) = \sum_{i \in T} Loss(\mathbf{x}_i^a, \mathbf{x}_i^p, \mathbf{x}_i^n)$$

set of all triplets



maximize
minimize

$f(\mathbf{x}^n)$

$f(\mathbf{x}^a)$

training

$f(\mathbf{x}^p)$

$f(\mathbf{x}^n)$

$f(\mathbf{x}^a)$

$f(\mathbf{x}^p)$

# Triplet Selection

- Motivation
  - Using all possible triples for training leads to slow convergence.
  - It selects hard triplets for training to ensure fast convergence.

$f(\mathbf{x}^n)$

$f(\mathbf{x}^p)$

$f(\mathbf{x}^a)$

**Easy Triplet**

$f(\mathbf{x}^n)$
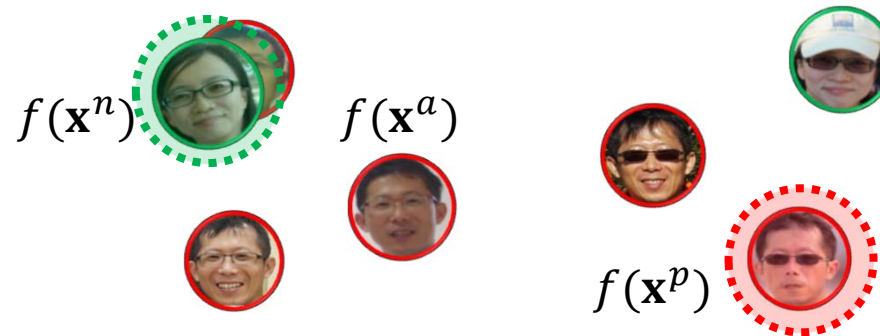
$f(\mathbf{x}^p)$

$f(\mathbf{x}^a)$

**Hard Triplet**

# Triplet Selection

- Intuitive Selection
  - find the hardest positive $\hat{\mathbf{x}}^p$ and negative $\hat{\mathbf{x}}^n$ to an anchor $\mathbf{x}^a$.

$$\hat{\mathbf{x}}^p = \arg\max_{\mathbf{x}^p} ||f(\mathbf{x}^a) - f(\mathbf{x}^p)|| \qquad \hat{\mathbf{x}}^n = \arg\min_{\mathbf{x}^n} ||f(\mathbf{x}^a) - f(\mathbf{x}^n)||$$



drawback: poor training
reason: poorly imaged faces would dominate the selection

# Triplet Selection

- Online Generation
  - mini-batch generation: sample training dataset to generate mini-batch around 1800 faces
    - around 40 faces per identity
    - randomly sampled negative faces
  - triplet formation: form all possible anchor-positive-negative triplets within the mini-batch.
    - use all anchor-positive pairs in mini-batch
    - randomly take a semi-hard negative (within margin)

# Triplet Selection



**Training Dataset** → **Sampling** → **Mini-Batch**

$\mathbf{x}^a$
all negatives e

$\mathbf{x}^p$
positive image

→ **FaceNet** →

$f(\mathbf{x}^p)$

$f(\mathbf{x}^a)$ $\alpha$

**semi-hard negative** $f(\mathbf{x}^n)$

20