



Boulder

Computer Vision; Face Recognition and Detection

Maziar Raissi

Assistant Professor

Department of Applied Mathematics

University of Colorado Boulder

maziar.raissi@colorado.edu



Boulder

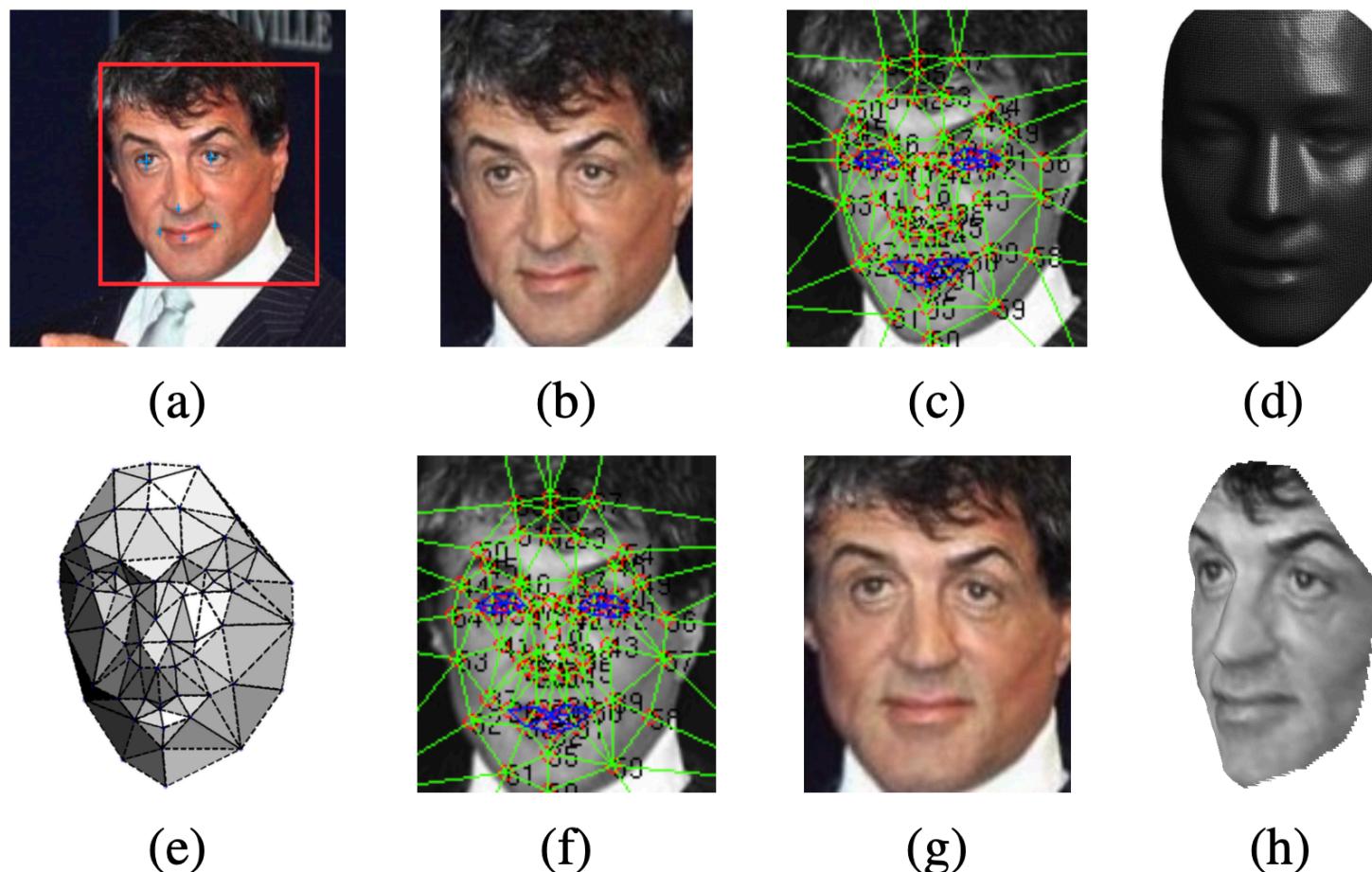
DeepFace: Closing the Gap to Human-Level Performance in Face Verification

detect \triangleright align \triangleright represent \triangleright classify

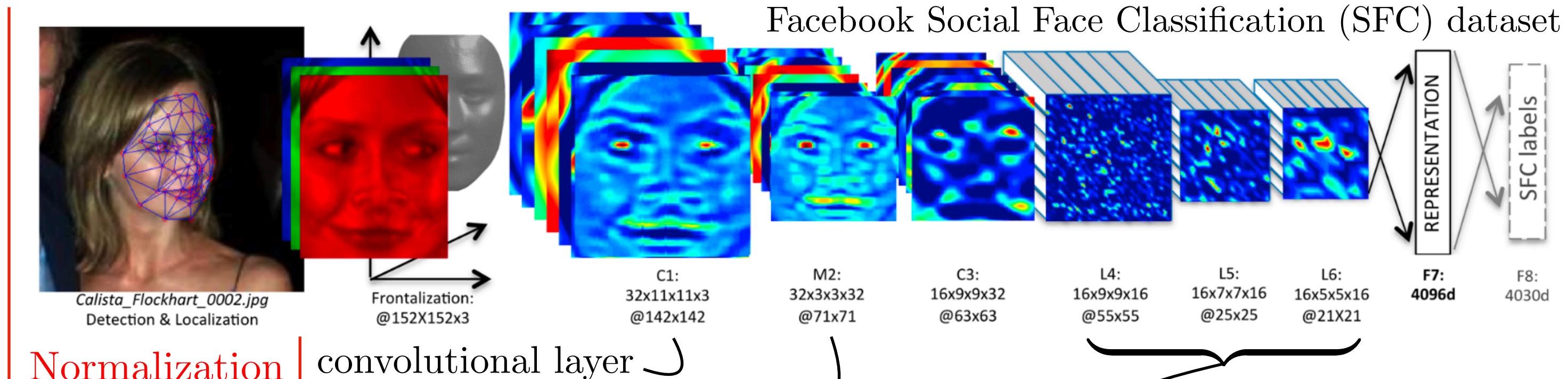
Model-based alignment: Explicit 3D face modeling in order to apply a piecewise affine transformation

Face representation from a nine-layer deep neural network

- Labeled Faces in the Wild (LFW) dataset
- YouTube Faces (YTF) dataset



Alignment pipeline. (a) The detected face, with 6 initial fiducial points. (b) The induced 2D-aligned crop. (c) 67 fiducial points on the 2D-aligned crop with their corresponding Delaunay triangulation, we added triangles on the contour to avoid discontinuities. (d) The reference 3D shape transformed to the 2D-aligned crop image-plane. (e) Triangle visibility w.r.t. to the fitted 3D-2D camera; darker triangles are less visible. (f) The 67 fiducial points induced by the 3D model that are used to direct the piece-wise affine warpping. (g) The final frontalized crop. (h) A new view generated by the 3D model (not used in this paper).



Normalization

$G(I) \rightarrow \text{representation}$

$$f(I) := \overline{G}(I)/\|\overline{G}(I)\|_2$$

$$\overline{G}(I)_i = \frac{G(I)_i}{\max(G(I)_i, \epsilon)}$$

convolutional layer

max-pooling layer

every location in the feature map learns a different set of filters

Since different regions of an aligned image have different local statistics, the spatial stationarity assumption of convolution cannot hold

Face verification metric: inner product between the two normalized feature vectors

$$\chi^2(f_1, f_2) = \sum_i w_i (f_1[i] - f_2[i])^2 / (f_1[i] + f_2[i])$$

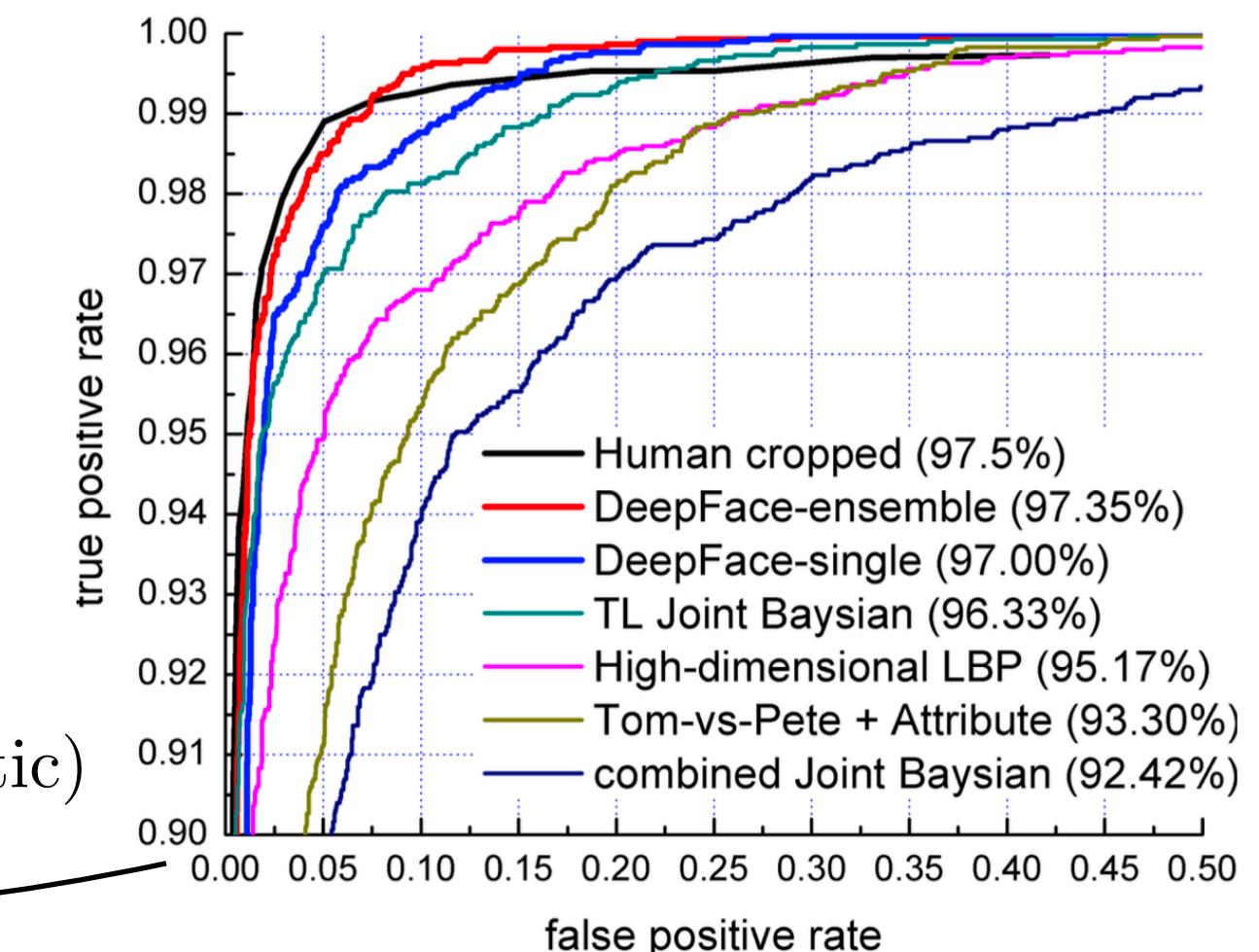
learned using a linear SVM

Weighted χ^2 distance

$$d(f_1, f_2) = \sum_i \alpha_i |f_1[i] - f_2[i]|$$

Siamese Network induced distance
(trained by cross entropy loss)

The ROC (Receiver Operating Characteristic)
curves on the LFW dataset





Boulder

FaceNet: A Unified Embedding for Face Recognition and Clustering

Face verification: Is this the same person?

Face recognition: Who is this person?

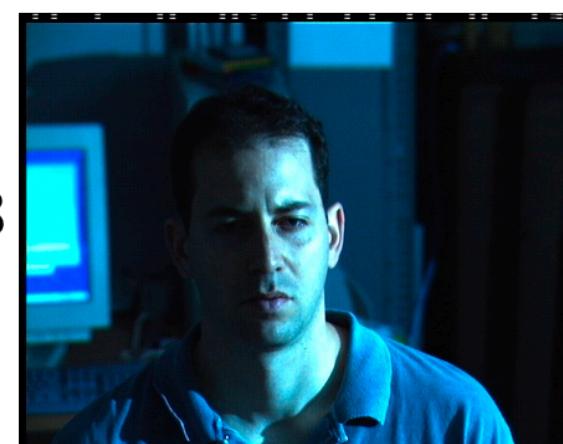
Clustering: Find common people among these faces?



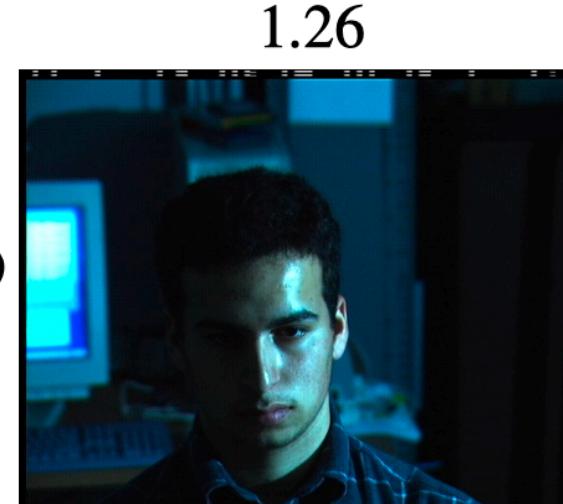
1.22



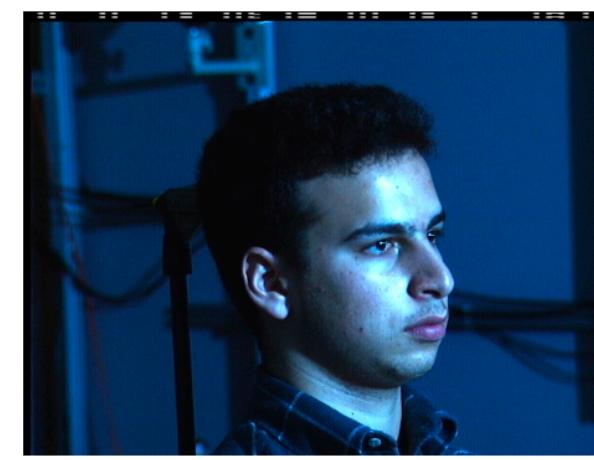
1.04



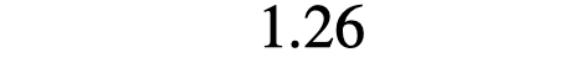
0.78



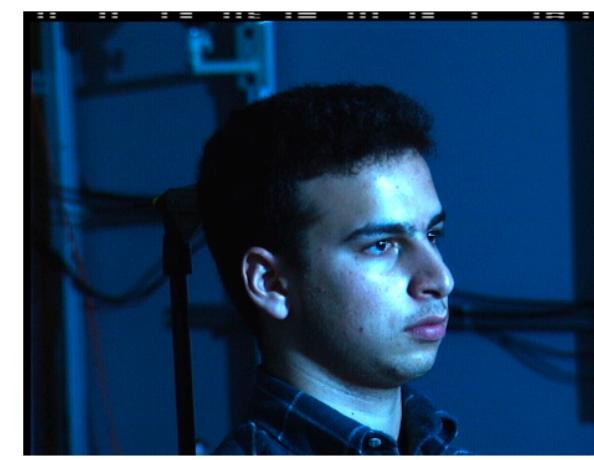
0.99



1.33



1.26



A threshold of 1.1 would classify every pair correctly



Triplet Loss

$x \mapsto f(x) \in \mathbb{R}^d$

embedding
image

$\|f(x)\|_2 = 1 \rightarrow$ constraint (hypersphere) (i.e., $f(x) = 0$)

- 1) The squared distance between all faces (independent of imaging conditions) of the same identity should be small
- 2) The squared distance between a pair of face images from different identities should be large

$x_i^a \rightarrow$ anchor (image of a specific person)

$x_i^p \rightarrow$ positive (other images of the same person)

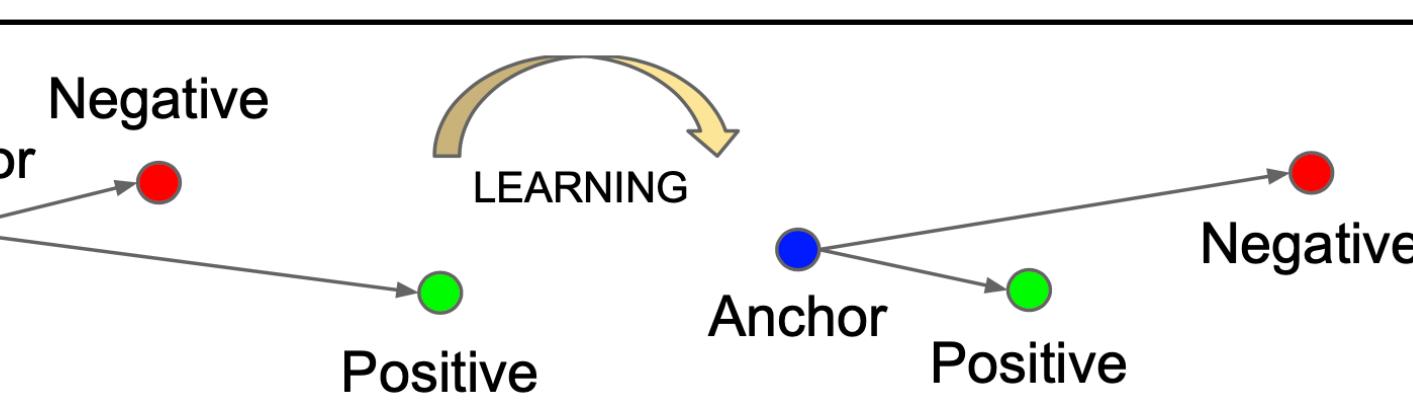
$x_i^n \rightarrow$ negative (any image of any other person)

$\|f(x_i^a) - f(x_i^p)\|_2^2 + \alpha < \|f(x_i^a) - f(x_i^n)\|_2^2, \forall (x_i^a, x_i^p, x_i^n) \in \mathcal{T}$

margin

set of all possible triplets

$$L = \sum_{i=1}^N [\|f(x_i^a) - f(x_i^p)\|_2^2 - \|f(x_i^a) - f(x_i^n)\|_2^2 + \alpha]_+$$



Triplet Selection

$$\arg \max_{x_i^p} \|f(x_i^a) - f(x_i^p)\|_2^2 \rightarrow \text{hard positive}$$

$$\arg \min_{x_i^n} \|f(x_i^a) - f(x_i^n)\|_2^2 \rightarrow \text{hard negative}$$

$$\|f(x_i^a) - f(x_i^p)\|_2^2 < \|f(x_i^a) - f(x_i^n)\|_2^2$$

semi-hard negative
avoid model collapse

Face Verification Task

$$D(x_i, x_j) \rightarrow L_2 \text{ distance threshold}$$

(same/different)

$\mathcal{P}_{\text{same}} \rightarrow$ all pairs (i, j) of the same identity

$\mathcal{P}_{\text{diff}} \rightarrow$ all pairs (i, j) of different identities

$\text{TA}(d) = \{(i, j) \in \mathcal{P}_{\text{same}}, \text{with } D(x_i, x_j) \leq d\}$

true accepts

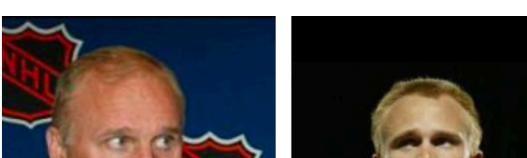
$\text{FA}(d) = \{(i, j) \in \mathcal{P}_{\text{diff}}, \text{with } D(x_i, x_j) \leq d\}$

false accepts

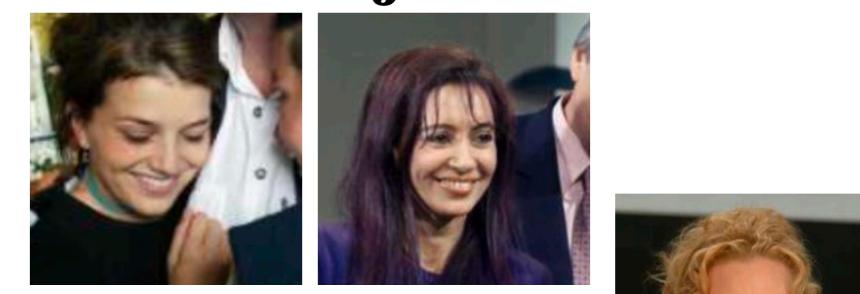
$$\text{VAL}(d) = \frac{|\text{TA}(d)|}{|\mathcal{P}_{\text{same}}|}, \quad \text{FAR}(d) = \frac{|\text{FA}(d)|}{|\mathcal{P}_{\text{diff}}|}$$

Validation rate and the false accept rate for a given face distance d

False accept



False reject



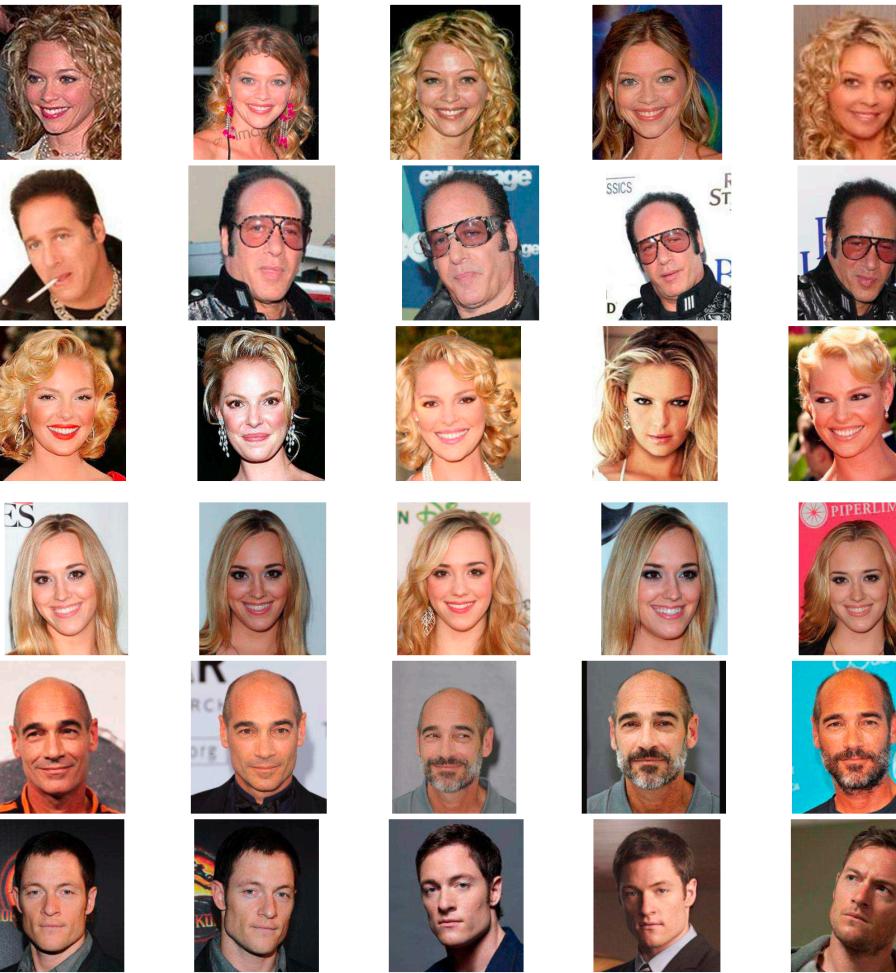


Boulder

Deep Face Recognition

Dataset	Identities	Images
LFW	5,749	13,233
WDRef [4]	2,995	99,773
CelebFaces [25]	10,177	202,599

Dataset	Identities	Images
Ours	2,622	2.6M
FaceBook [29]	4,030	4.4M
Google [17]	8M	200M



Stage	Aim	Type	# of persons	# of images per person	total # images	Annotation effort	100%-EER
1	Candidate list generation	A	5,000	200	1,000,000	—	—
2	Image set expansion	M	2,622	2,000	5,244,000	4 days	—
3	Rank image sets	A	2,622	1,000	2,622,000	—	96.90
4	Near dup. removal	A	2,622	623	1,635,159	—	—
5	Final manual filtering	M	2,622	375	982,803	10 days	92.83

Dataset Collection

- IMDB (Internet Movie Data Base)
- Freebase Knowledge Graph
- Google Image Search

No.	Config	Data	Train Align.	Test Align.	Embedding	100% - EER
1	A	C	No	No	No	92.83
2	A	F	No	No	No	95.80
3	A	F	No	Yes	No	96.70
4	B	F	No	Yes	No	97.27
5	B	F	Yes	Yes	No	96.17
6	D	F	No	Yes	No	96.73
7	B	F	No	Yes	Yes	99.13

The top 50 images (based on Google search rank in the downloaded set) for each identity are used as positive training samples, and the top 50 images of all other identities are used as negative training samples. A one-vs-rest linear SVM is trained for each identity using the Fisher Vector Faces descriptor.

VLAD descriptor ▷ Clustering ▷ Near Duplicate Removal

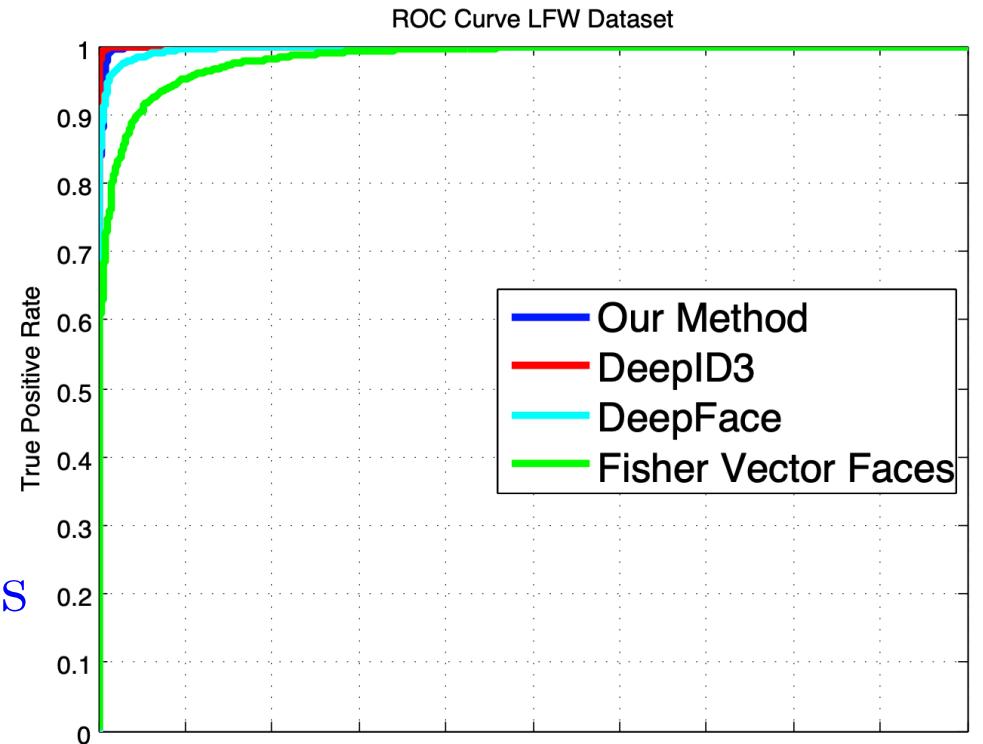
Final manual filtering is aided with AlexNet trained to discriminate between the 2,622 face identities.

layer type name	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
input	conv	relu	conv	relu	mpool	conv	relu	conv	relu	mpool	conv	relu	conv	relu	conv	relu	mpool	conv	
–	conv1_1	relu1_1	conv1_2	relu1_2	pool1	conv2_1	relu2_1	conv2_2	relu2_2	pool2	conv3_1	relu3_1	conv3_2	relu3_2	conv3_3	relu3_3	pool3	conv4_1	
support	–	3	1	3	1	2	3	1	3	1	2	3	1	3	1	2	3	1	
filt dim	–	3	–	64	–	–	64	–	128	–	–	128	–	256	–	256	–	256	
num filters	–	64	–	64	–	–	128	–	128	–	–	256	–	256	–	256	–	512	
stride	–	1	1	1	1	2	1	1	1	1	2	1	1	1	1	1	2	1	
pad	–	1	0	1	0	0	1	0	1	0	0	1	0	1	0	1	0	1	
layer type name	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37
relu	conv	relu	conv	relu	mpool	conv	relu	conv	relu	conv	relu	mpool	conv	relu	conv	relu	conv	softmax	
relu4_1	conv4_2	relu4_2	conv4_3	relu4_3	pool4	conv5_1	relu5_1	conv5_2	relu5_2	conv5_3	relu5_3	pool5	fc6	relu6	fc7	relu7	fc8	prob	
support	1	3	1	3	1	2	3	1	3	1	2	7	1	1	1	1	1	1	
filt dim	–	512	–	512	–	–	512	–	512	–	–	512	–	4096	–	4096	–	–	
num filters	–	512	–	512	–	–	512	–	512	–	–	4096	–	4096	–	2622	–	–	
stride	1	1	1	1	1	2	1	1	1	1	2	1	1	1	1	1	1	1	
pad	0	1	0	1	0	0	1	0	1	0	0	0	0	0	0	0	0	0	

Youtube Faces Dataset

EER: Equal Error Rate
error rate @ the ROC operating point where false positive & false negative rates are equal

No.	Method	Images	Networks	Acc.
1	Fisher Vector Faces [21]	-	-	93.10
2	DeepFace [29]	4M	3	97.35
3	Fusion [30]	500M	5	98.37
4	DeepID-2,3		200	99.47
5	FaceNet [17]	200M	1	98.87
6	FaceNet [17] + Alignment	200M	1	99.63
7	Ours	2.6M	1	98.95



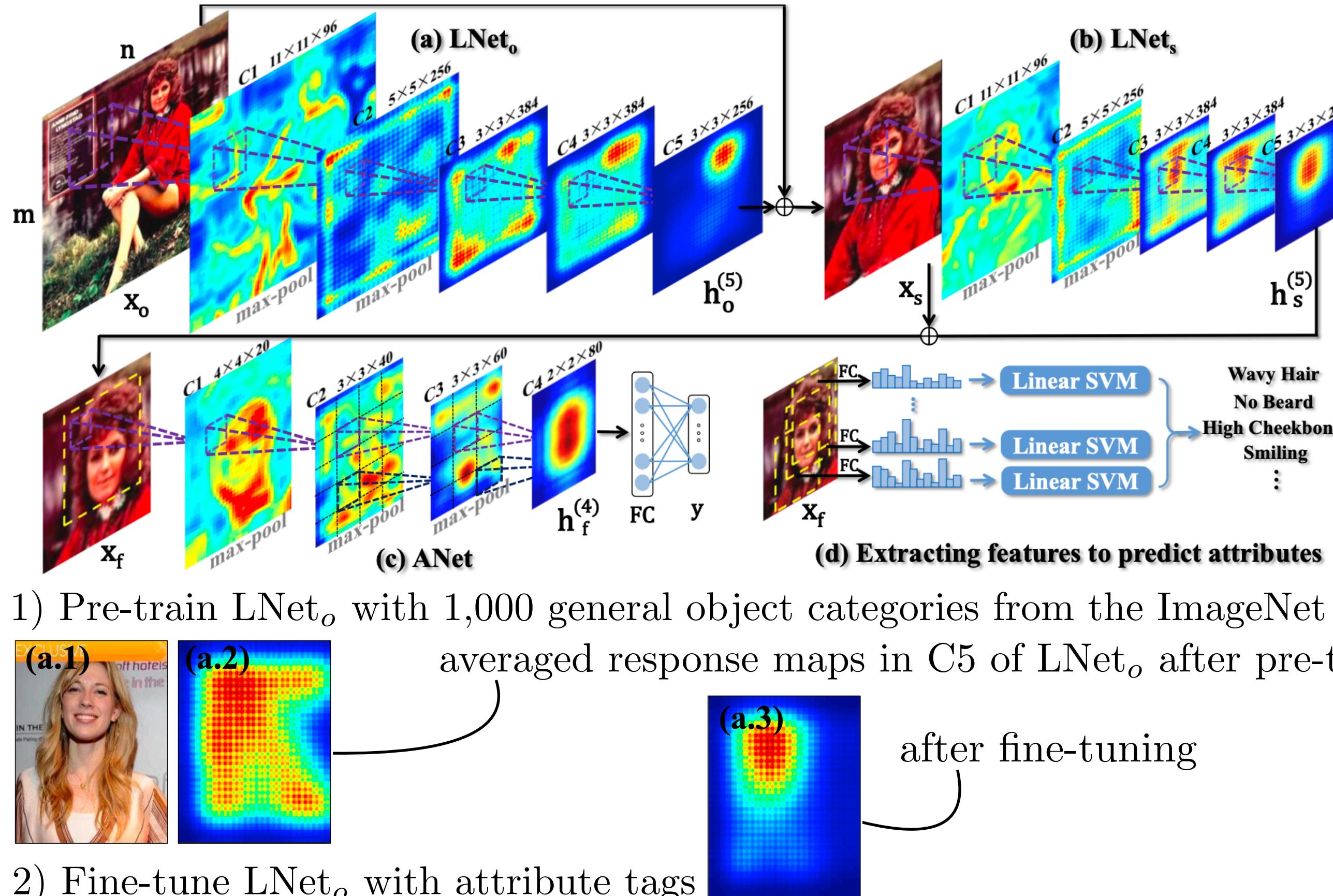
K indicates the number of faces used to represent each video

No.	Method	Images	Networks	100%- EER	Acc.
1	Video Fisher Vector Faces [15]	-	-	87.7	83.8
2	DeepFace [29]	4M	1	91.4	91.4
3	DeepID-2,2+,3		200	-	93.2
4	FaceNet [17] + Alignment	200M	1	-	95.1
5	Ours (<i>K</i> = 100)	2.6M	1	92.8	91.6
6	Ours (<i>K</i> = 100) + Embedding learning	2.6M	1	97.4	97.3



Boulder

Deep Learning Face Attributes in the Wild



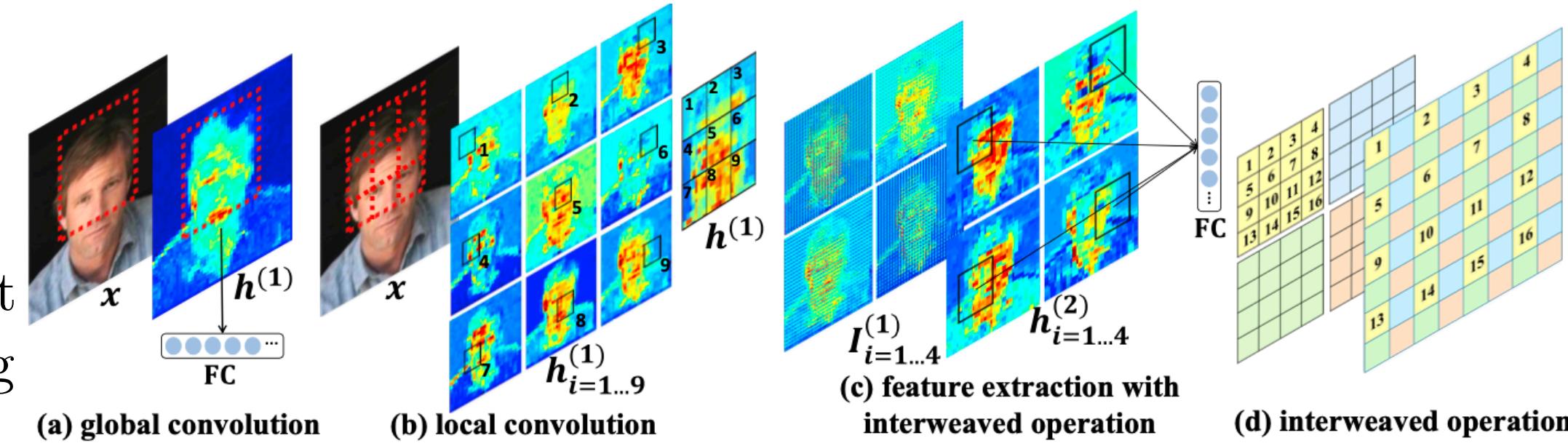
4) Repeat steps 1-3 for $LNet_s$

5) Pre-train ANet with a large number of face identities from the CelebFaces dataset

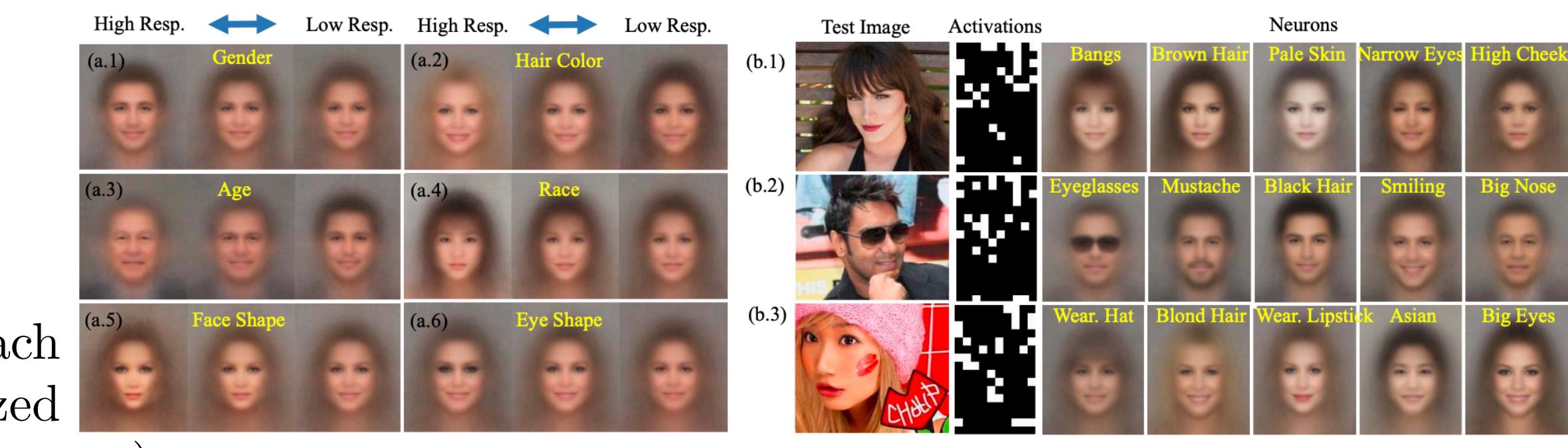
ANet is pre-trained by combining the softmax loss and the similarity loss

$$\mathcal{L} = \sum_{i,y_i=y_j} \|FC_i - FC_j\|_2^2 \rightarrow \text{improve intra-class invariance}$$

6) ANet is fine-tuned by attributes to learn high-level feature FC



Detailed pipeline of efficient feature extractions in ANet.

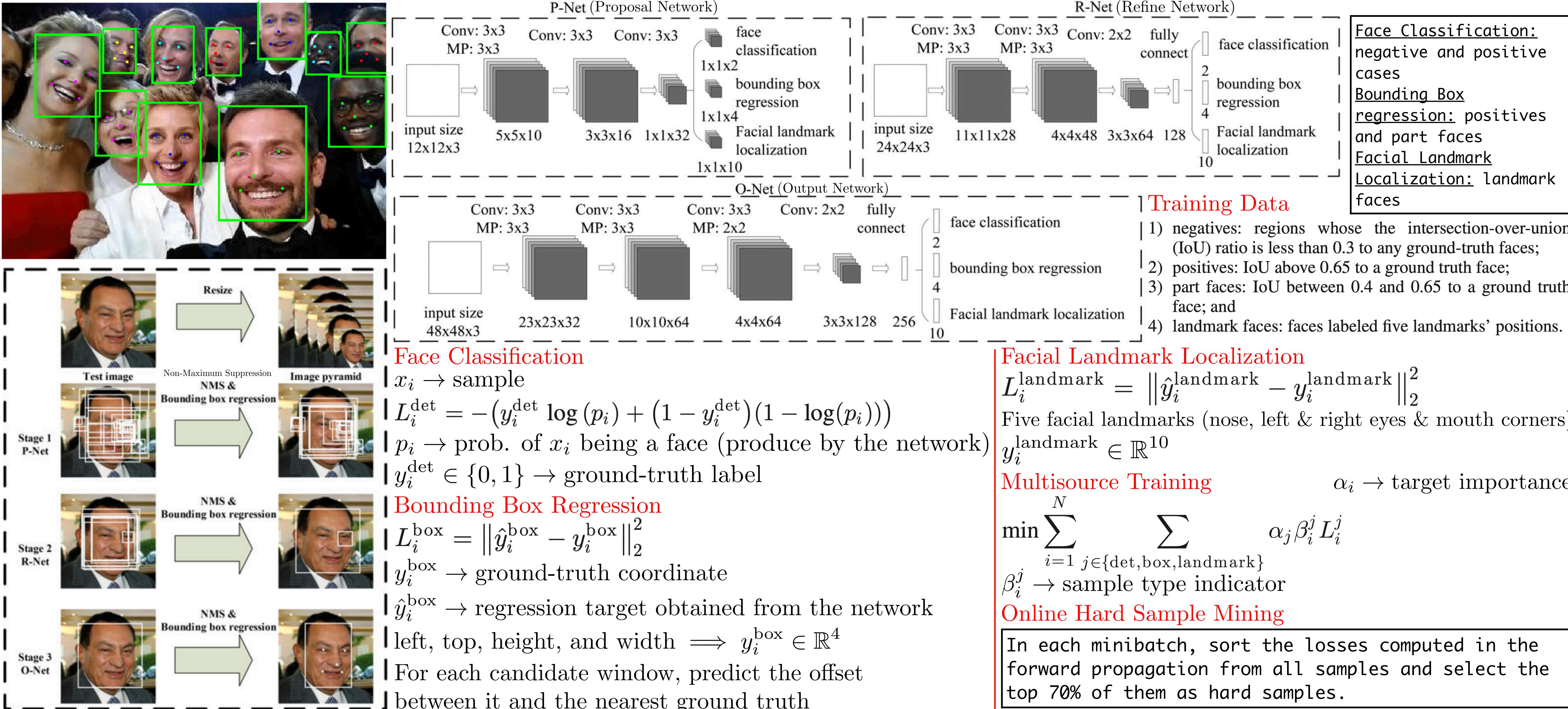


Visualization of neurons in ANet (a) after pre-training (b) after fine-tuning



Boulder

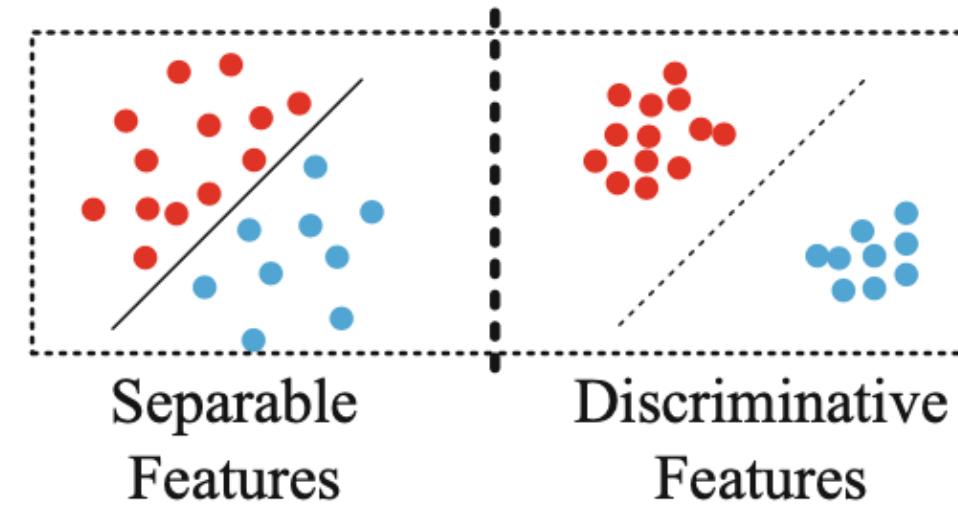
Joint Face Detection and Alignment Using Multitask Cascaded Convolutional Networks





Boulder

A Discriminative Feature Learning Approach for Deep Face Recognition



- inter-class dispersion
- intra-class compactness

$$\mathcal{L}_S = - \sum_{i=1}^m \log \frac{e^{W_{y_i}^T \mathbf{x}_i + b_{y_i}}}{\sum_{j=1}^n e^{W_j^T \mathbf{x}_i + b_j}} \rightarrow \text{softmax loss}$$

$\mathbf{x}_i \in \mathbb{R}^d \rightarrow$ deep features of i -th example belonging to class y_i

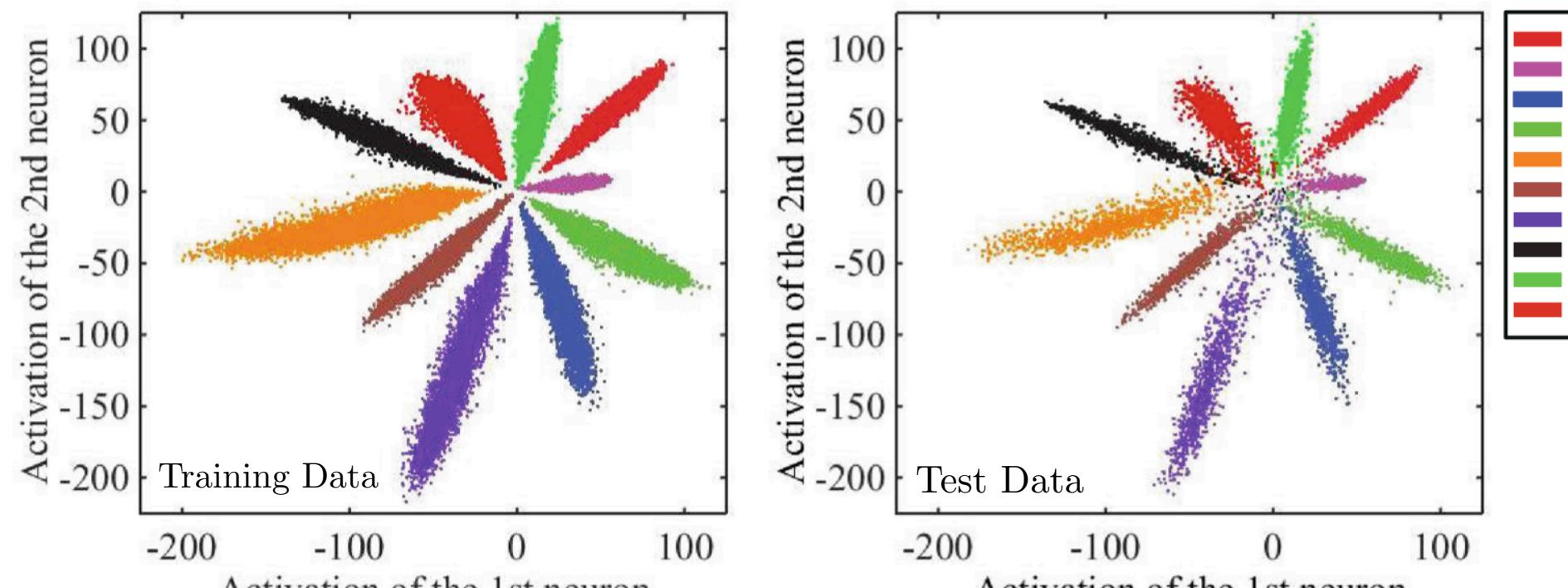
$W_j \in \mathbb{R}^d \rightarrow j$ -th column of the weights $W \in \mathbb{R}^{d \times n}$

$b \in \mathbb{R}^n \rightarrow$ bias

$m \rightarrow$ size of mini-batch

$n \rightarrow$ number of classes

MNIST Example (CNN : $I_i \mapsto \mathbf{x}_i \in \mathbb{R}^2$)



- Labeled Faces in the Wild (LFW)
- YouTube Faces (YTF)
- MegaFace Challenge

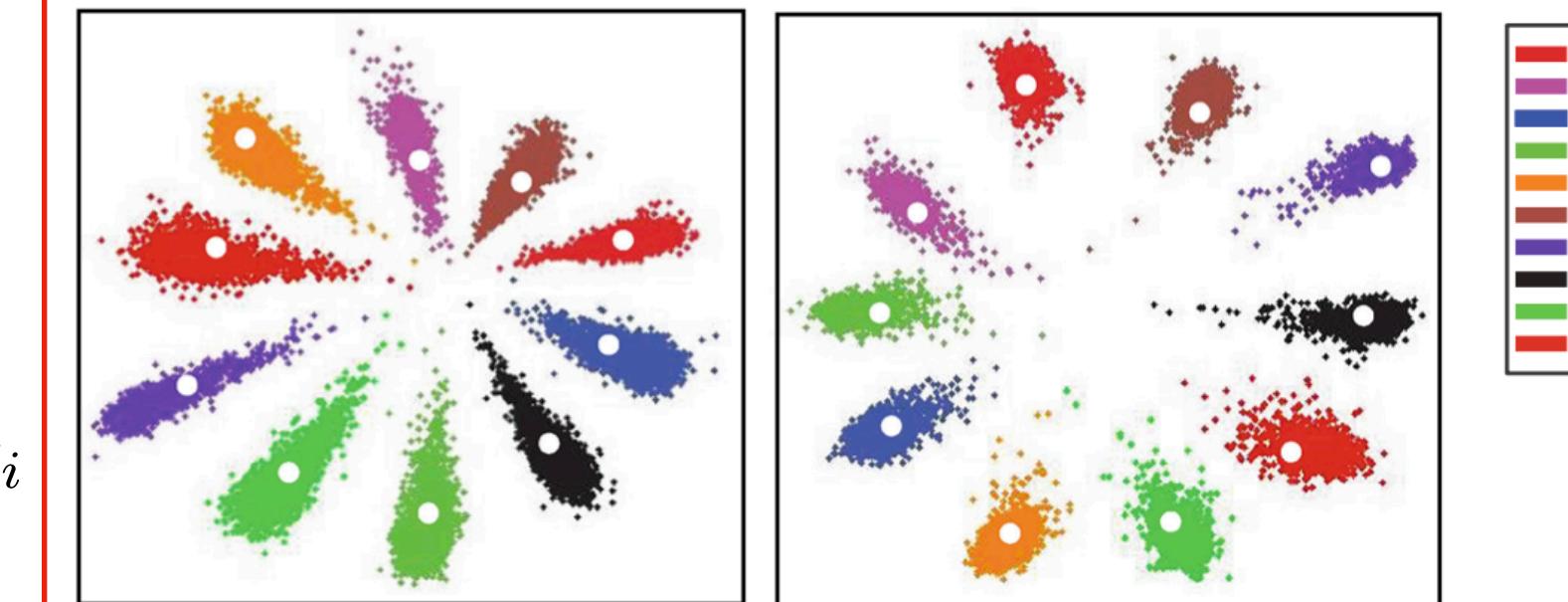
- Face Recognition
- Face Verification

The Center Loss

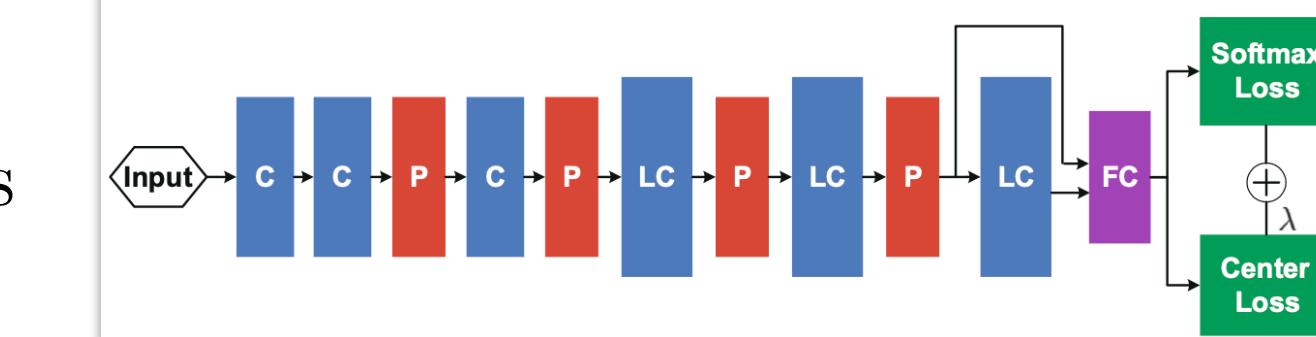
$$\mathcal{L}_C = \frac{1}{2} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{c}_{y_i}\|_2^2$$

$\mathbf{c}_{y_i} \in \mathbb{R}^d \rightarrow y_i$ th class center of deep features

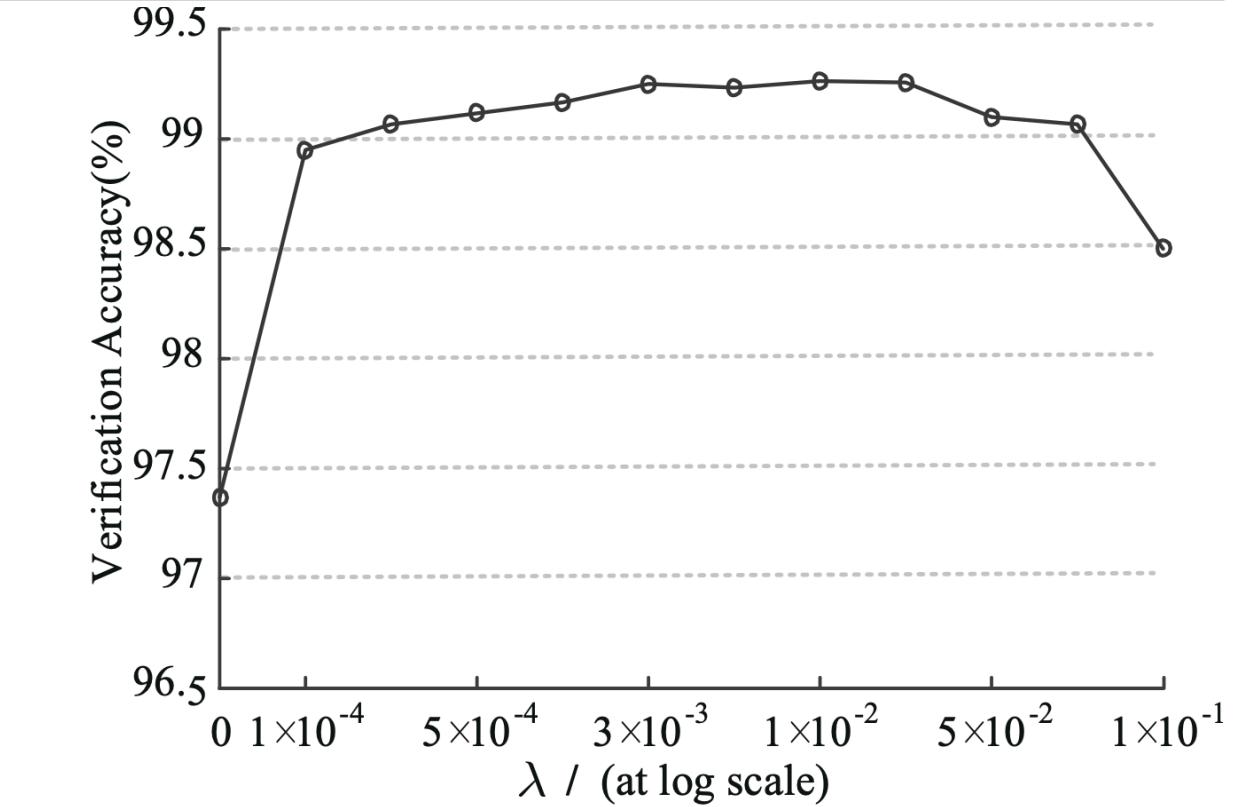
$$\mathcal{L} = \mathcal{L}_S + \lambda \mathcal{L}_C$$



C: The convolution layer
P: The max-pooling layer
LC: The local convolution layer
FC: The fully connected layer



Method	Images	Networks	Acc. on LFW	Acc. on YTF
DeepFace [34]	4M	3	97.35 %	91.4 %
DeepID-2+ [32]	-	1	98.70 %	-
DeepID-2+ [32]	-	25	99.47 %	93.2 %
FaceNet [27]	200M	1	99.63 %	95.1 %
Deep FR [25]	2.6M	1	98.95 %	97.3 %
Baidu [21]	1.3M	1	99.13 %	-
model A	0.7M	1	97.37 %	91.1 %
model B	0.7M	1	99.10 %	93.8 %
model C (Proposed)	0.7M	1	99.28 %	94.9 %





Boulder

ArcFace: Additive Angular Margin Loss for Deep Face Recognition

- softmax classifier
- triplet loss

softmax drawbacks

- size of the linear transformation matrix $W \in \mathbb{R}^{d \times n}$ increases linearly with the identities number n
- good for closed-set classification problems but not discriminative enough for the open-set face recognition problem

triplet loss drawbacks

- combinatorial explosion in the number of face triplets
- semi-hard-example-mining is a quite difficult problem for efficient model training

$$L_1 = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{W_{y_i}^T x_i + b_{y_i}}}{\sum_{j=1}^n e^{W_j^T x_i + b_j}}$$

softmax loss

$x_i \in \mathbb{R}^d \rightarrow$ deep feature of the i -th sample

$y_i \rightarrow$ class of the i -th sample

$W_j \in \mathbb{R}^d \rightarrow j$ -th column of $W \in \mathbb{R}^{d \times n}$

$b_j \in \mathbb{R}^n \rightarrow$ bias

N & $n \rightarrow$ batch size & class number

set $b_j = 0$

$$W_j^T x_i = \|W_j\| \|x_i\| \cos \theta_j$$

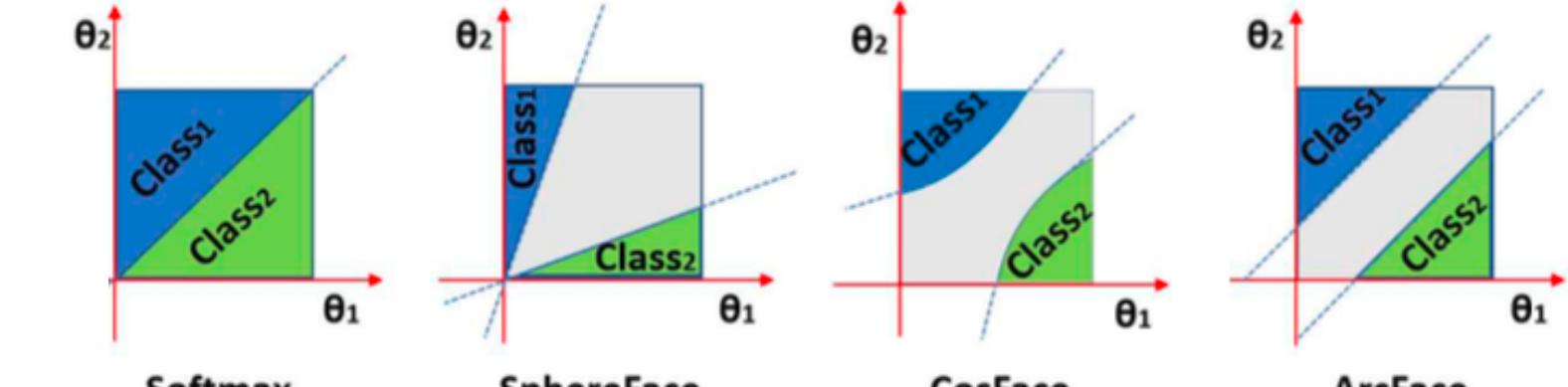
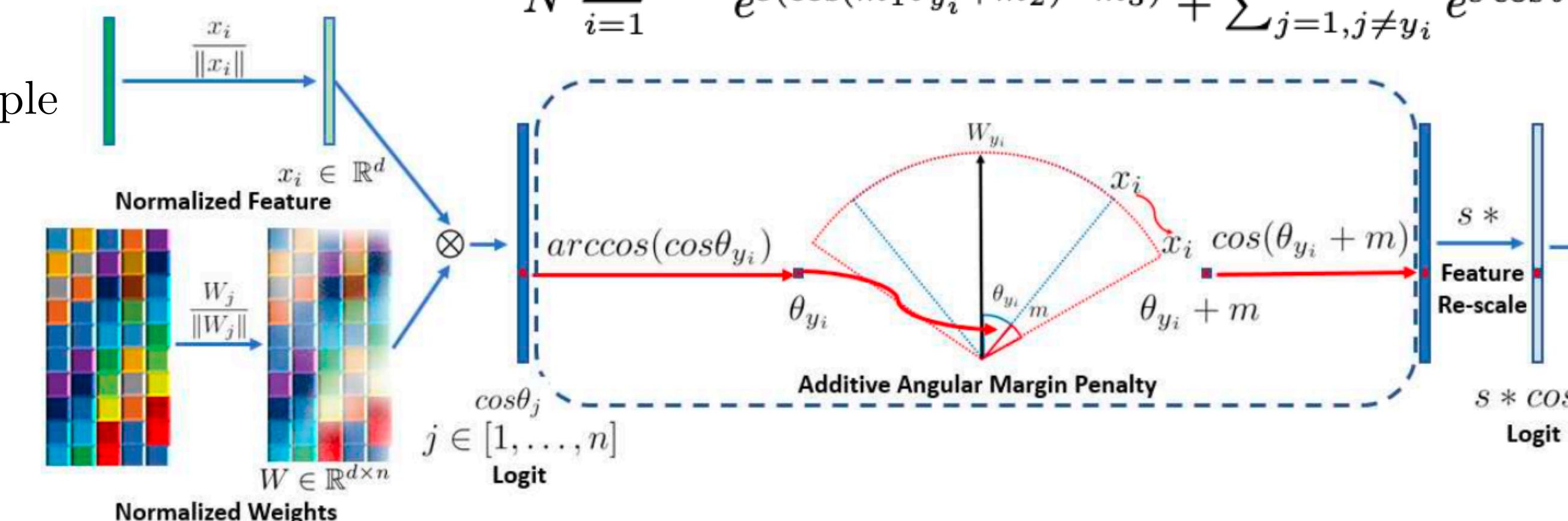
$\|W_j\| = 1$ by l_2 normalization
 $\|x_i\| = s$ by l_2 normalization
 ↪ hyper-sphere with radius s

$$L_2 = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s \cos \theta_{y_i}}}{e^{s \cos \theta_{y_i}} + \sum_{j=1, j \neq y_i}^n e^{s \cos \theta_j}}$$

$$L_3 = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s(\cos(\theta_{y_i} + m))}}{e^{s(\cos(\theta_{y_i} + m))} + \sum_{j=1, j \neq y_i}^n e^{s \cos \theta_j}}$$

$m \rightarrow$ additive angular margin penalty between x_i and W_{y_i} to simultaneously enhance the intra-class compactness and inter-class discrepancy

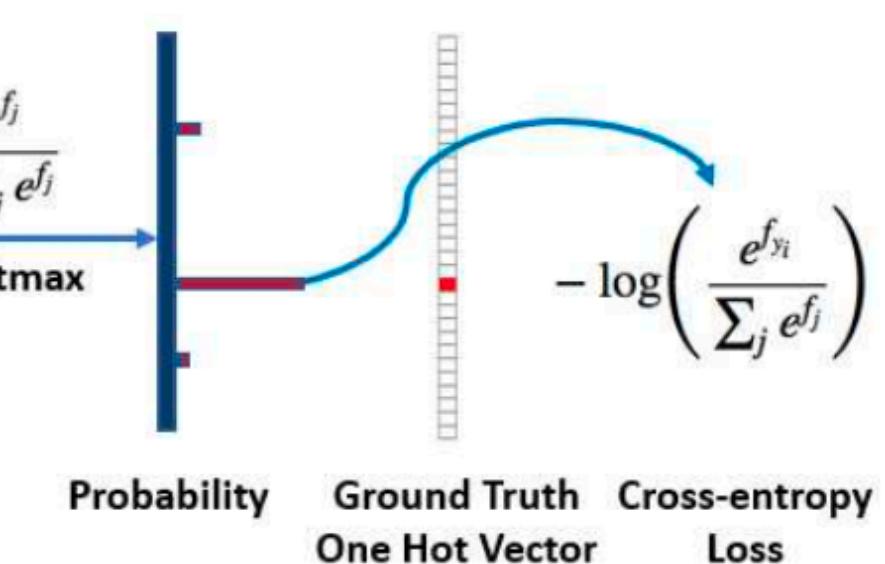
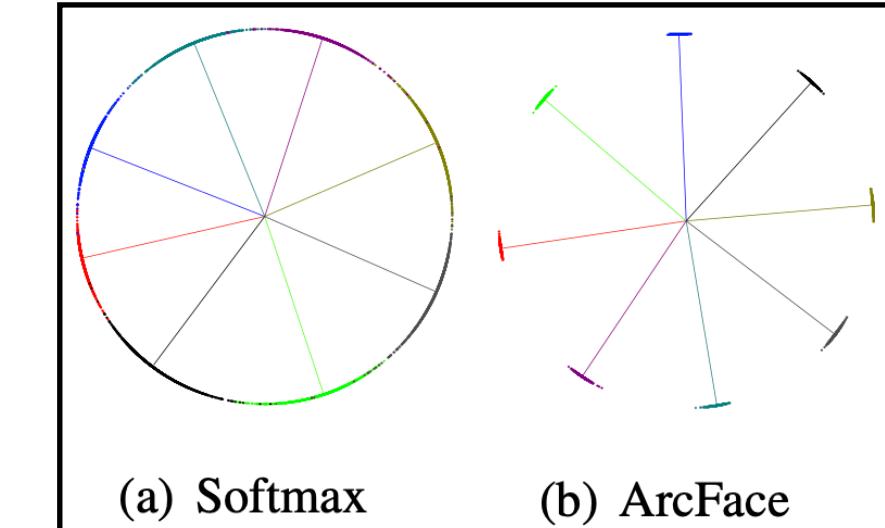
$$L_4 = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s(\cos(m_1 \theta_{y_i} + m_2) - m_3)}}{e^{s(\cos(m_1 \theta_{y_i} + m_2) - m_3)} + \sum_{j=1, j \neq y_i}^n e^{s \cos \theta_j}}$$



triplet-loss $\leftarrow \arccos(x_i^{pos} x_i) + m \leq \arccos(x_i^{neg} x_i)$

Intra-Loss $\leftarrow L_5 = L_2 + \frac{1}{\pi N} \sum_{i=1}^N \theta_{y_i}$

Inter-Loss $\leftarrow L_6 = L_2 - \frac{1}{\pi N(n-1)} \sum_{i=1}^N \sum_{j=1, j \neq y_i}^n \arccos(W_{y_i}^T W_j)$





Boulder

Questions?
