



香港中文大學

The Chinese University of Hong Kong

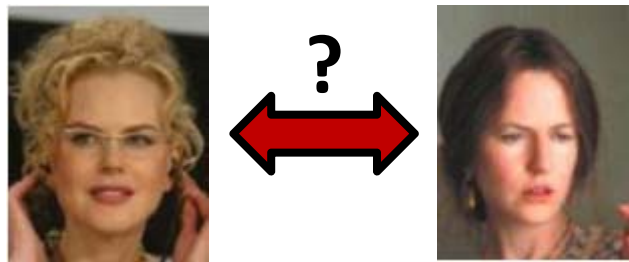
DeepID: Deep Learning for Face Recognition

Xiaogang Wang

Department of Electronic Engineering,
The Chinese University of Hong Kong

Face Recognition

- Face verification: binary classification
 - Verify two images belonging to the same person or not



- Face identification: multi-class classification
 - classify an image into one of N identity classes



Labeled Faces in the Wild (2007)



Random guess (50%)

Best results
without deep learning

MSRA TL Joint Bayesian (96.33%)

Human funneled (99.20%)

CUHK deep learning result (99.53%)
Google deep learning result (99.6%)



Deep Learning Results on LFW

Method	Accuracy (%)	# points	# training images
Huang et al. CVPR'12	87%	3	Unsupervised
Sun et al. ICCV'13	92.52%	5	87,628
Facebook (CVPR'14)	97.35%	6 + 67	7,000,000
DeepID (CVPR'14)	97.45%	5	202,599
DeepID2 (NIPS'14)	99.15%	18	202,599
DeepID2+ (CVPR'15)	99.47%	18	450,000
Google (CVPR'15)	99.63%		200,000,000

- The first deep learning work on face recognition was done by Huang et al. in 2012. With unsupervised learning, the accuracy was 87%
- Our work at ICCV'13 achieved result (92.52%) comparable with state-of-the-art
- Our work at CVPR'14 reached **97.45%** close to “human cropped” performance (**97.53%**)
- DeepFace developed by Facebook also at CVPR'14 used 73-point 3D face alignment and 7 million training data (35 times larger than us)
- Our NIPS'14 work reached **99.15%** close to “human funneled” performance (**99.20%**)

Closed- and open-set face identification on LFW

Method	Rank-1 (%)	DIR @ 1% FAR (%)
COST-S1 [1]	56.7	25
COST-S1+s2 [1]	66.5	35
DeepFace [2]	64.9	44.5
DeepFace+ [3]	82.5	61.9
DeepID2	91.1	61.6
DeepID2+	95.0	80.7

[1] L. Best-Rowden, H. Han, C. Otto, B. Klare, and A. K. Jain. Unconstrained face recognition: Identifying a person of interest from a media collection. *TR MSU-CSE-14-1*, 2014.

[2] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. DeepFace: Closing the gap to human-level performance in face verification. In *Proc. CVPR*, 2014.

[3] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Web-scale training for face identification. Technical report, arXiv:1406.5266, 2014.

Learn face representations from

face verification, identification, multi-view reconstruction

Properties of face representations

sparseness, selectiveness, robustness

Applications of face representations

face localization, attribute recognition

Learn face representations from

face verification, identification, multi-view reconstruction

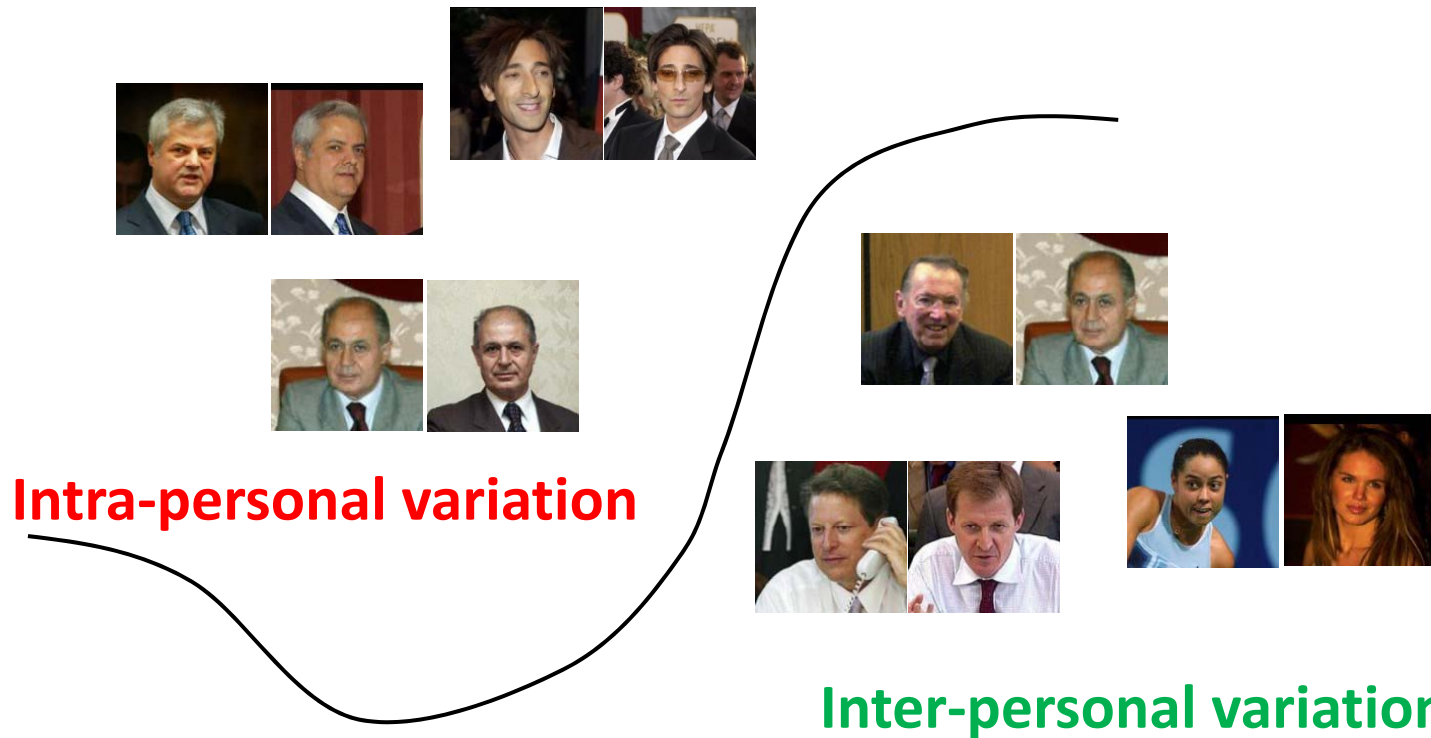
Properties of face representations

sparseness, selectiveness, robustness

Applications of face representations

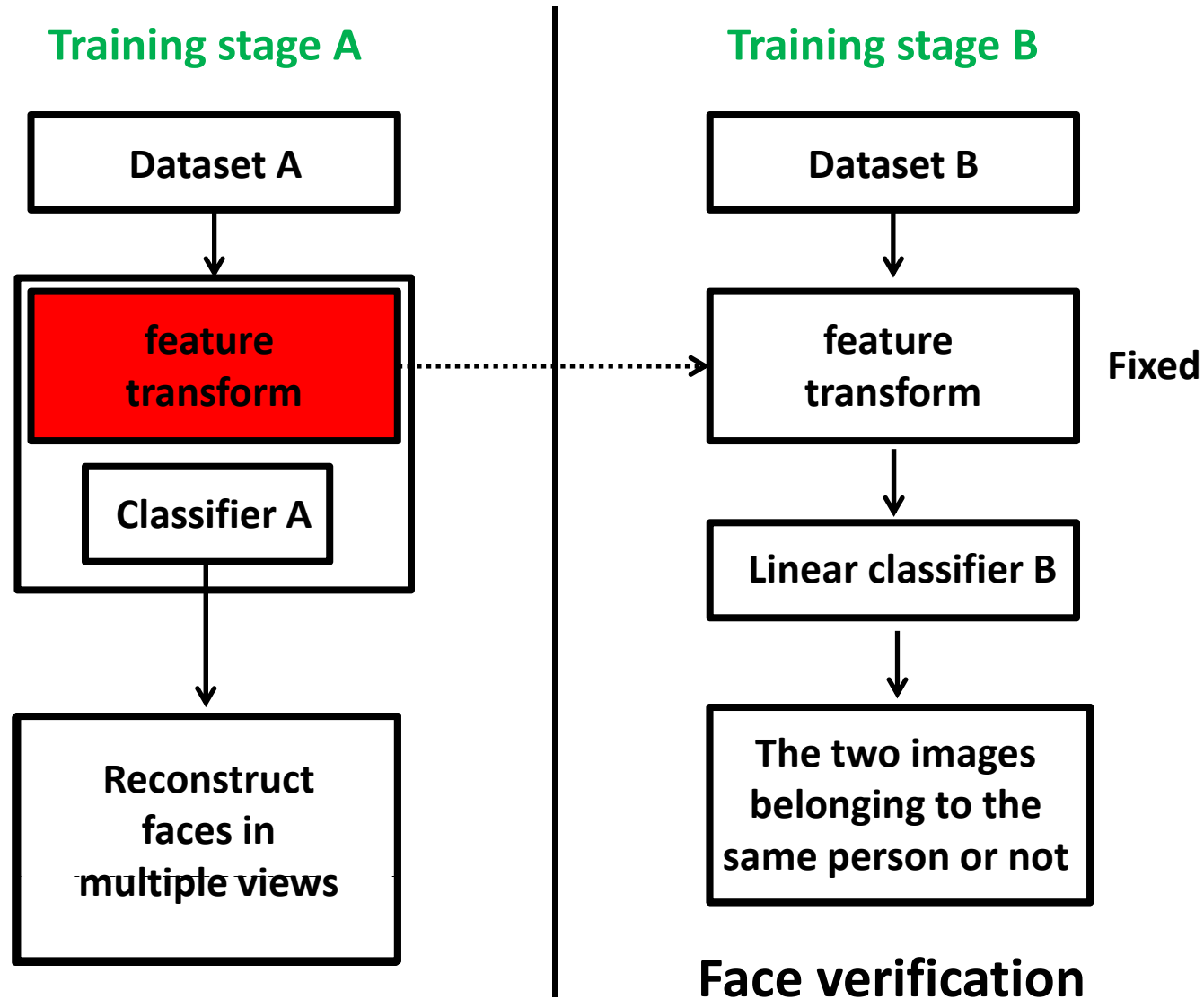
face localization, attribute recognition

Key challenge on face recognition



How to separate the two types of variations?

Learning feature representations



Learn face representations from

Prediction becomes richer

*Prediction becomes more
challenging*

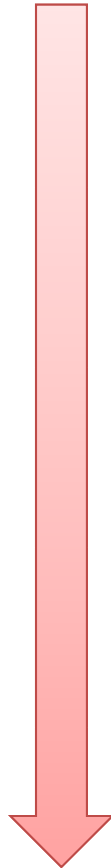
Supervision becomes stronger

*Feature learning becomes
more effective*

Predicting binary labels (verification)

Predicting multi-class labels (identification)

**Predicting thousands of real-valued pixels
(multi-view) reconstruction**

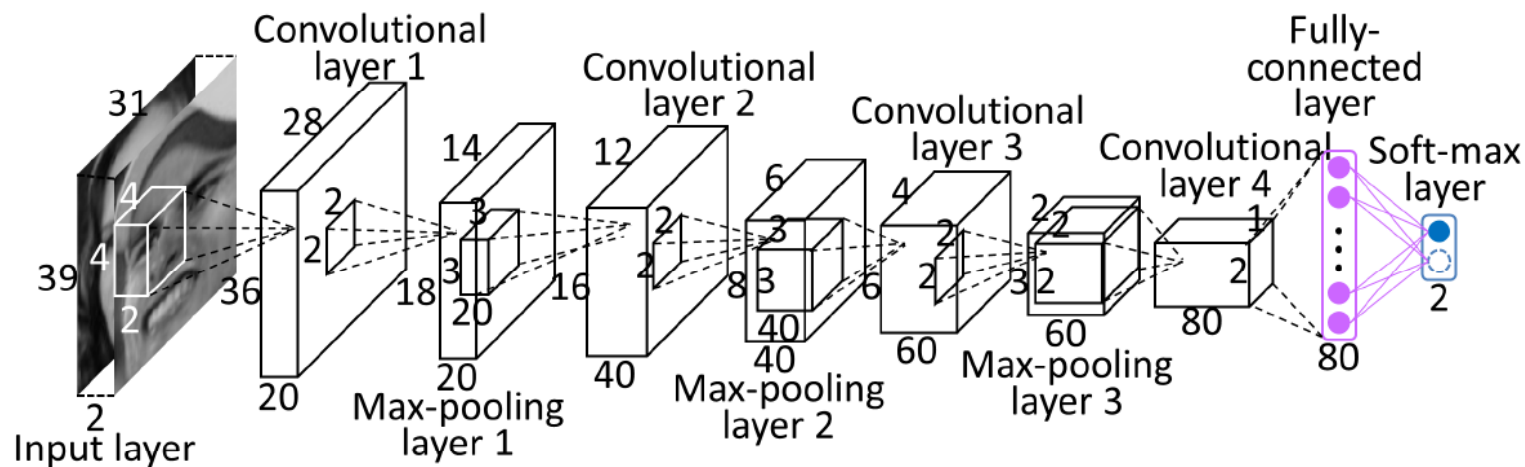


Learn face representations with verification signal

- Extract relational features with learned filter pairs

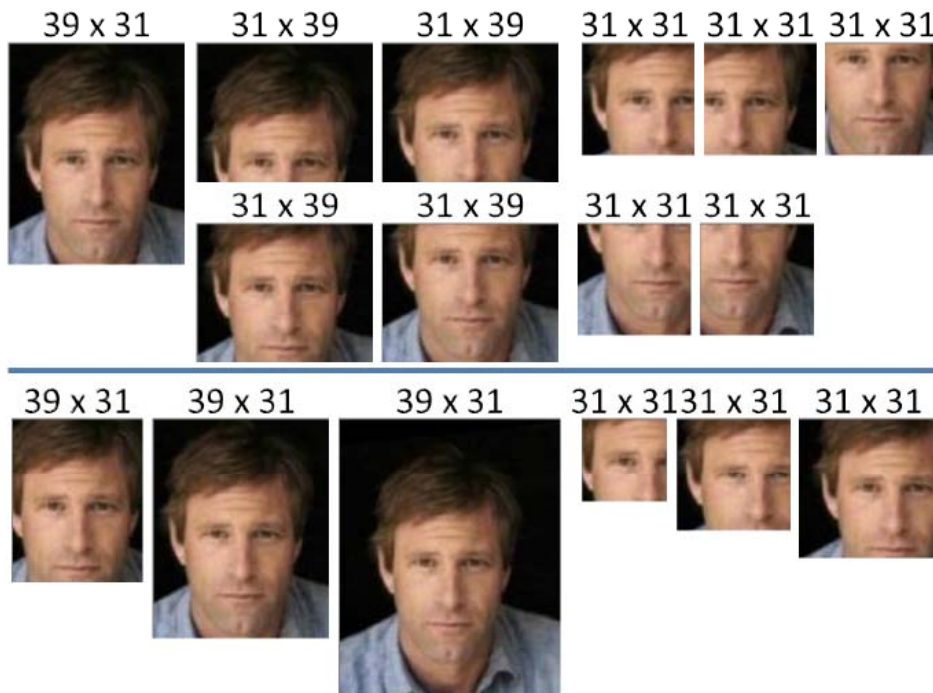
$$y^j = f(b^j + k^{1j} * x^1 + k^{2j} * x^2)$$

- These relational features are further processed through multiple layers to extract global features
- The fully connected layer is the feature representation

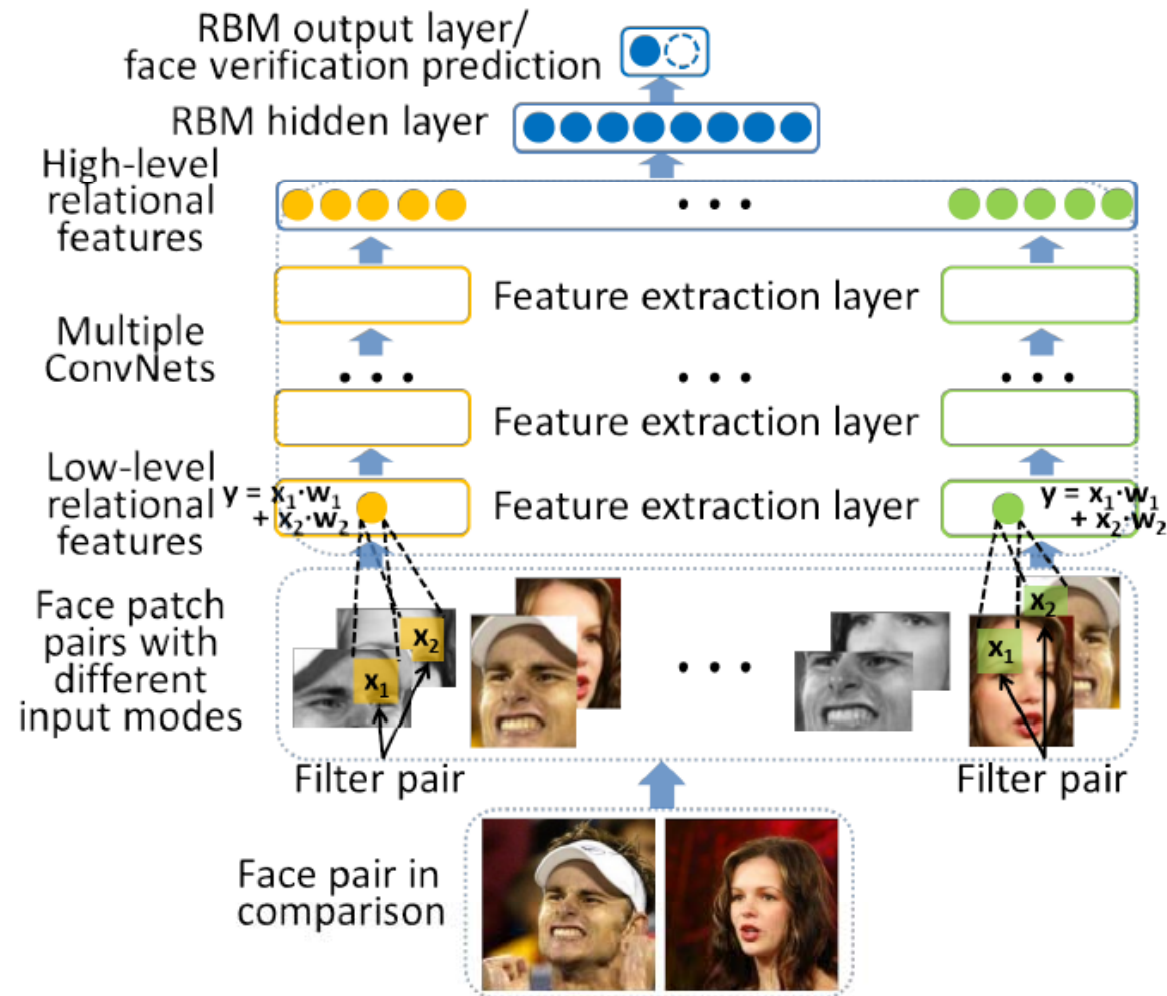


Generate multiple CNNs

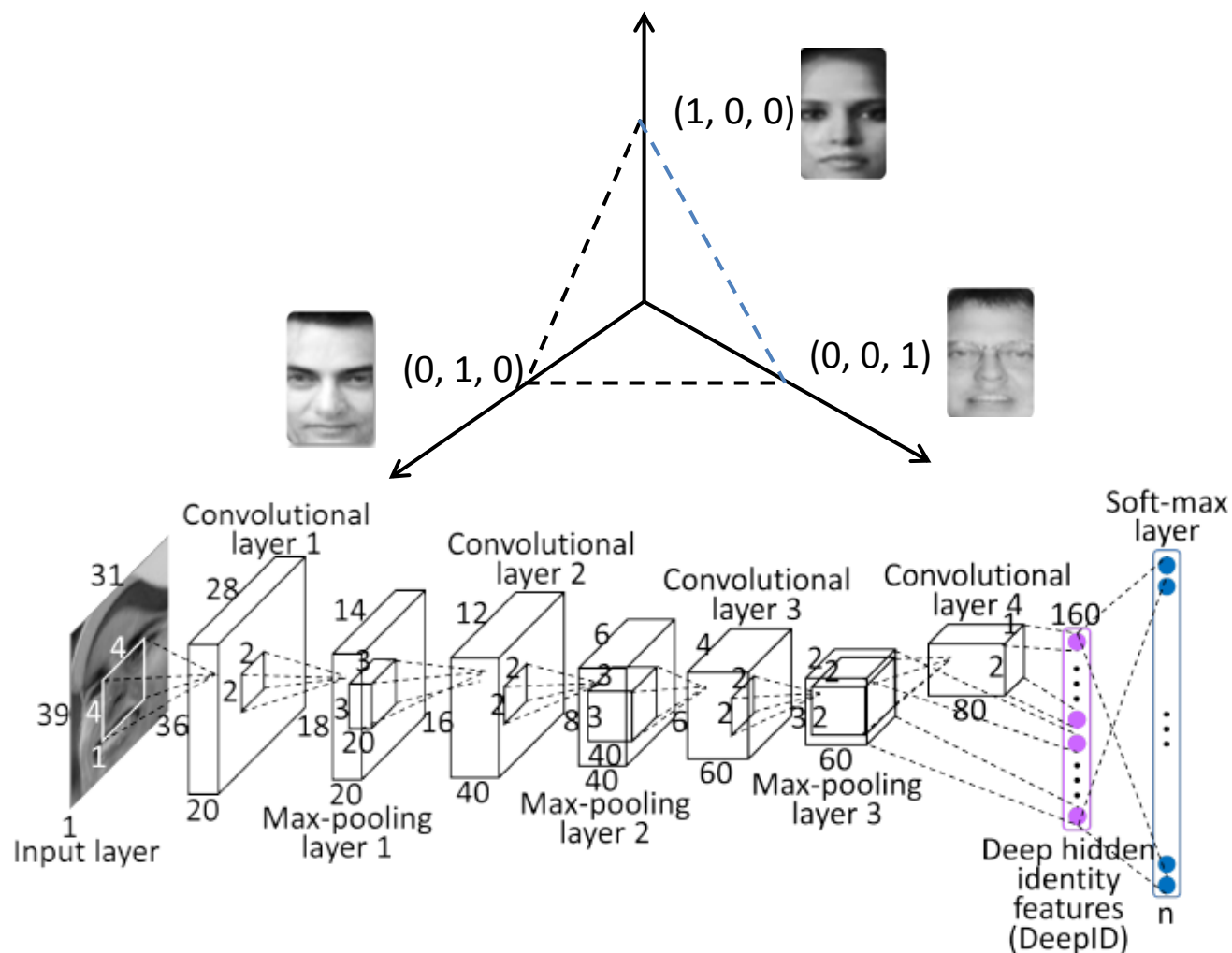
- 10 face regions, 3 scales, color/gray and 8 modes
- Base on three-point alignment



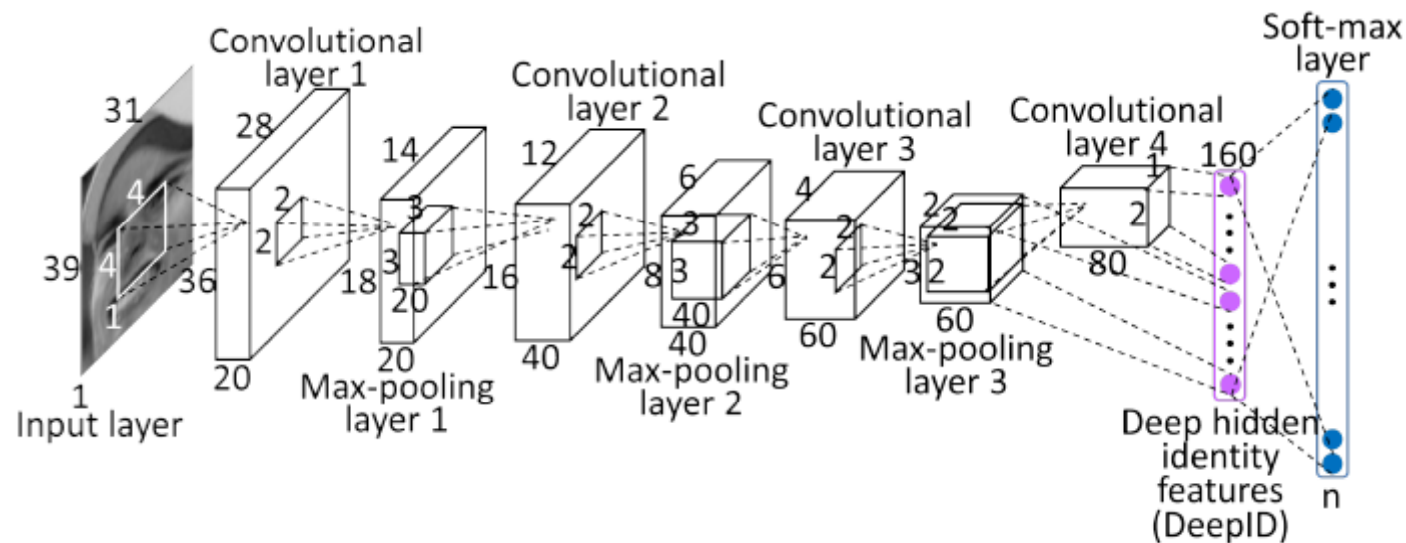
RBM combines features extracted by multiple CNNs



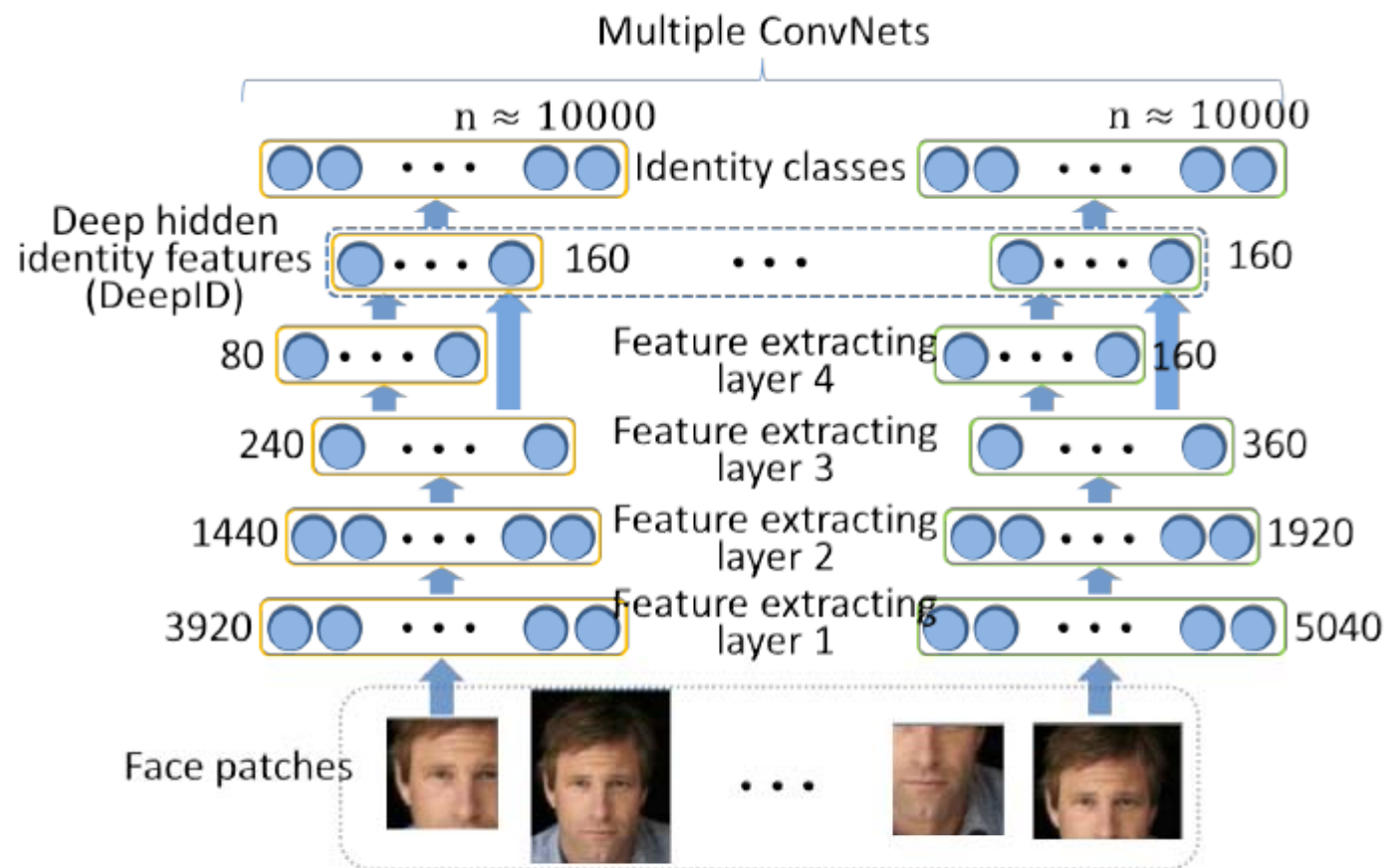
DeepID: Learn face representations with identification signal



- Features are from the last two convolution layers
- Learned features keep rich inter-personal variations
- Features can be well generalized to other tasks (e.g. verification) and identities outside the training set
- Increasing the number of classes to be predicted, the generalization power of the learned features improves



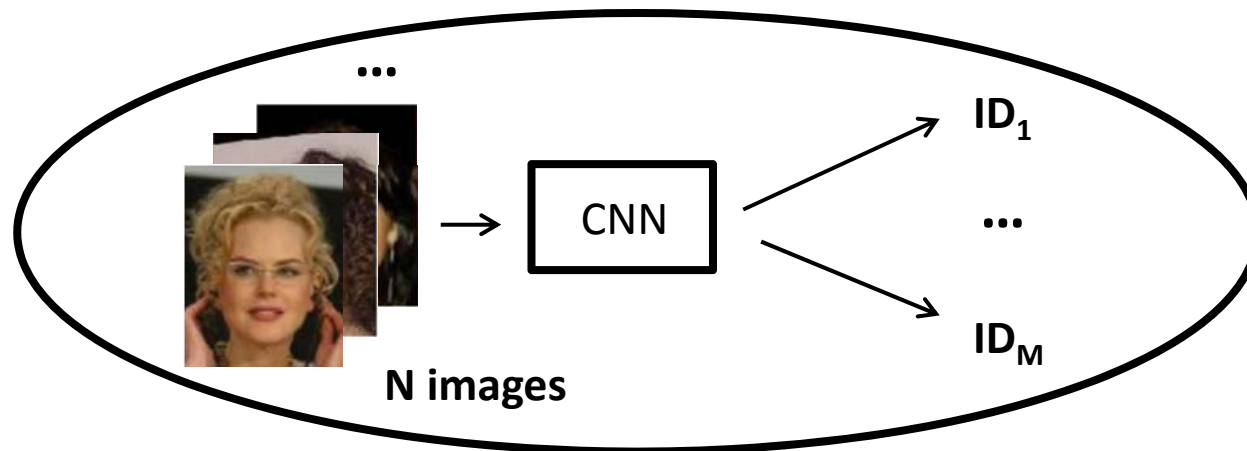
Extract features from multiple ConvNets



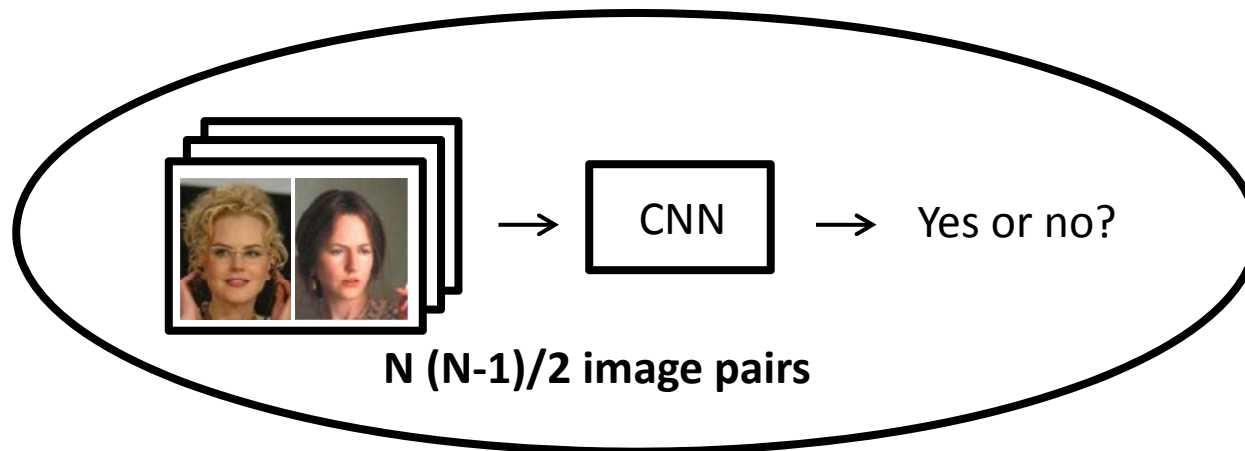
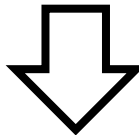
Learn feature representations with identification signal

- These features can be further processed by other classifiers in face verification. Interestingly, we find *Joint Bayesian* is more effective than cascading another neural network to classify these features

Why using identification as supervision is more efficiency than verification?



Same amount of labeling information



Identification supervision:

effective on capturing inter-personal variation

Verification supervision:

effective on reducing intra-personal variation

DeepID2: Joint identification-verification signals

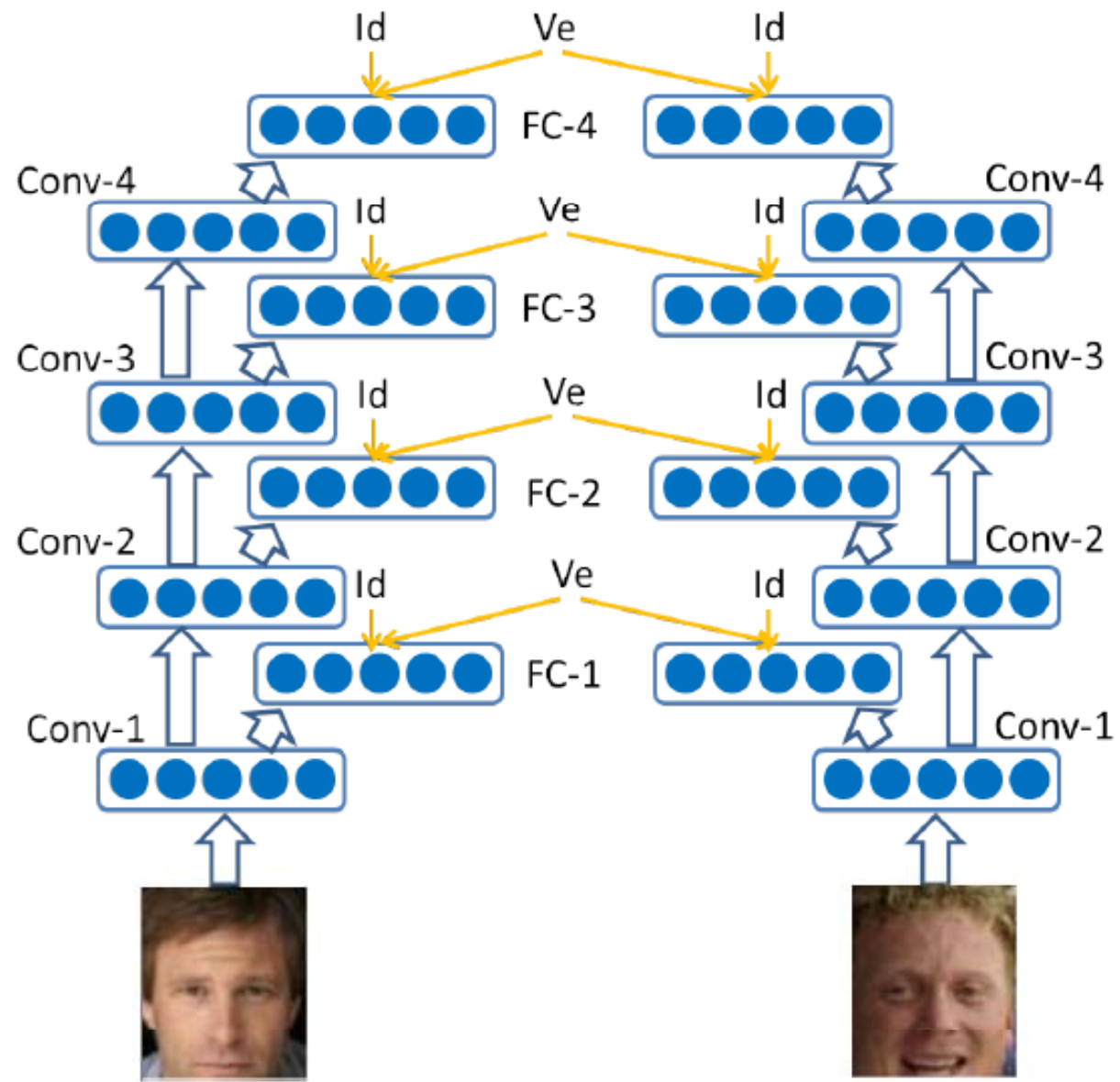
- Every two feature vectors extracted from the same identity should be close to each other

$$\text{Verif}(f_i, f_j, y_{ij}, \theta_{ve}) = \begin{cases} \frac{1}{2} \|f_i - f_j\|_2^2 & \text{if } y_{ij} = 1 \\ \frac{1}{2} \max(0, m - \|f_i - f_j\|_2)^2 & \text{if } y_{ij} = -1 \end{cases}$$

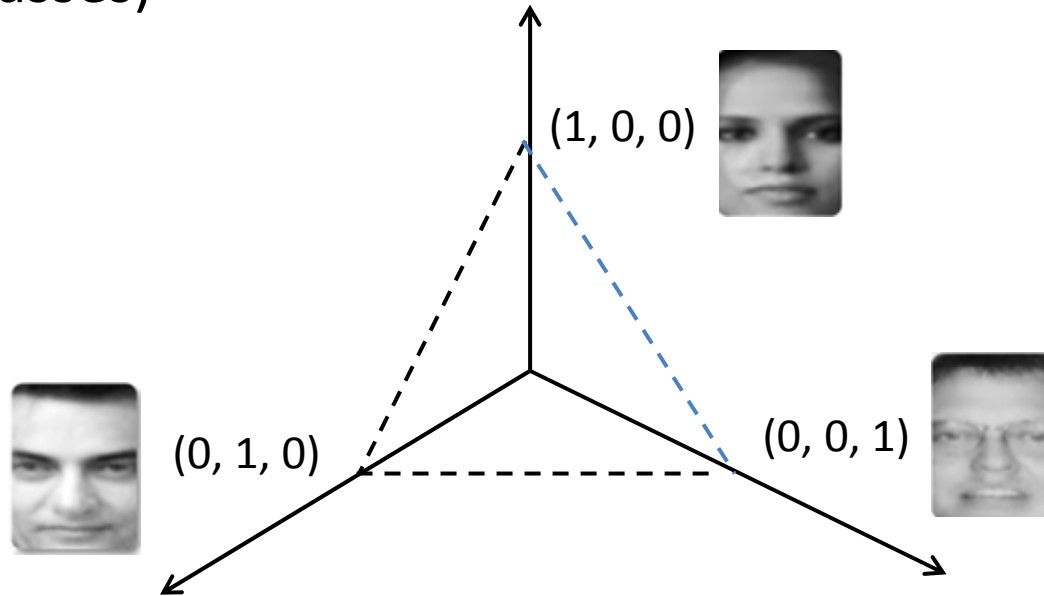
f_i and f_j are feature vectors extracted from two face images in comparison

$y_{ij} = 1$ means they are from the same identity; $y_{ij} = -1$ means different identities

m is a margin to be learned



Minimize the intra-personal variation under the constraint that the distance between classes is constant (i.e. contracting the volume of the image space without reducing the distance between classes)

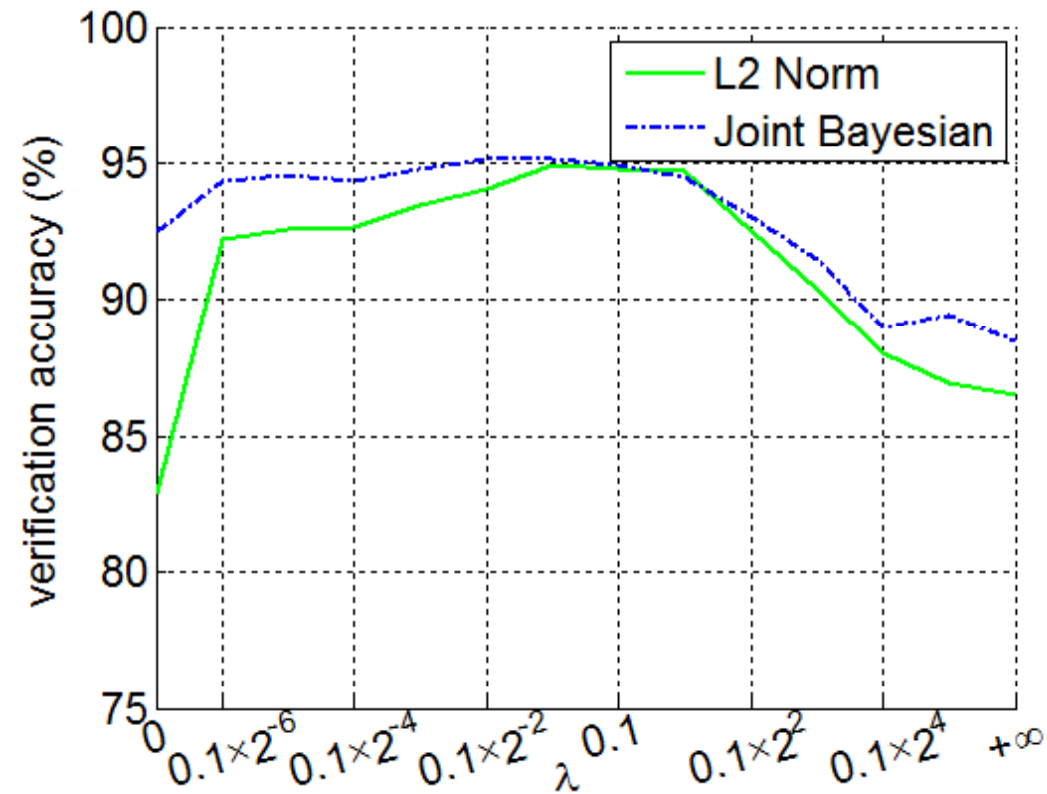


$$\mathbf{y} = f(\mathbf{x}); \quad g = \text{softmax}()$$

$$f^* = \arg \min_f \sum_{(i,j) \in \Omega_I} \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|^2$$

$$\text{s.t. } |g(f(\mathbf{x}_i)) - g(f(\mathbf{x}_j))| = 1, \quad \text{label}(\mathbf{x}_i) \neq \text{label}(\mathbf{x}_j)$$

Balancing identification and verification signals with parameter λ

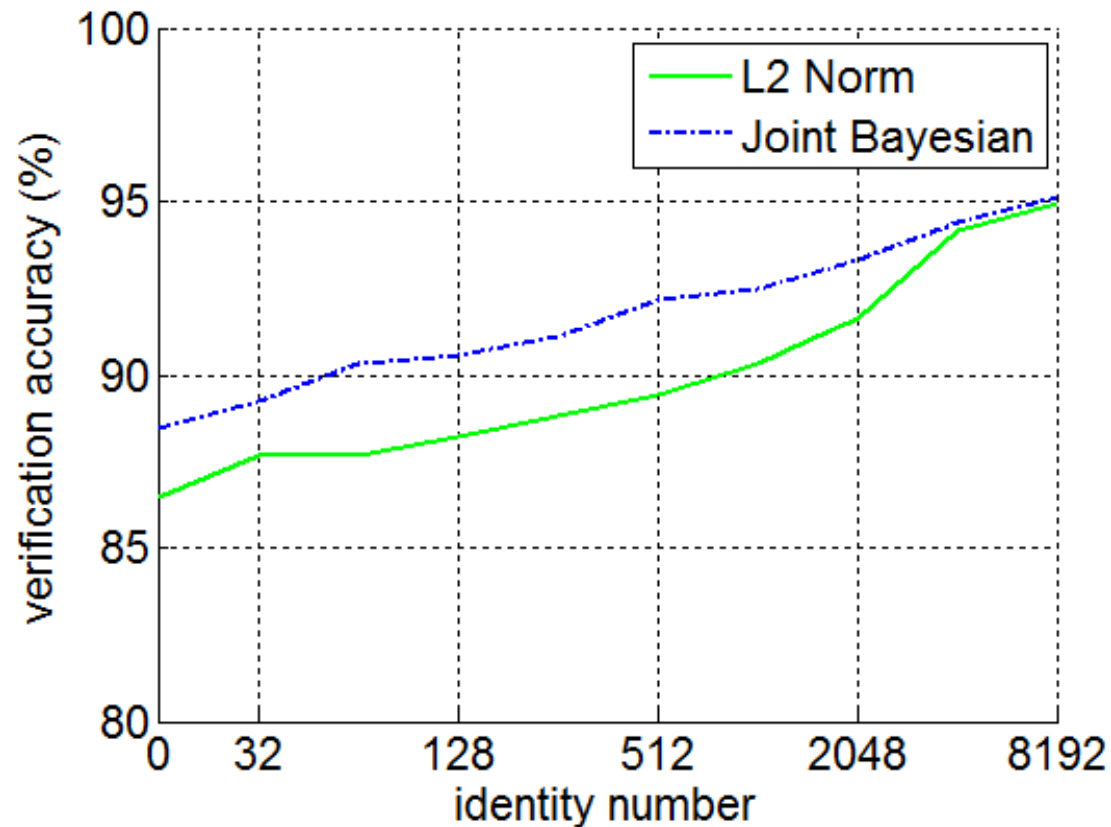


$\lambda = 0$: only identification signal

$\lambda = +\infty$: only verification signal

Rich identity information improves feature learning

- Face verification accuracies with the number of training identities



Summary of DeepID2

- 25 face regions at different scales and locations around landmarks are selected to build 25 neural networks
- All the 160×25 hidden identity features are further compressed into a 180-dimensional feature vector with PCA as a signature for each image
- With a single Titan GPU, the feature extraction process takes 35ms per image

Final Result on LFW

Methods	High-dim LBP [1]	TL Joint Bayesian [2]	DeepFace [3]	DeepID [4]	DeepID2 [5]	DeepID2+
Accuracy (%)	95.17	96.33	97.35	97.45	99.15	99.47

[1] Chen, Cao, Wen, and Sun. Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification. CVPR, 2013.

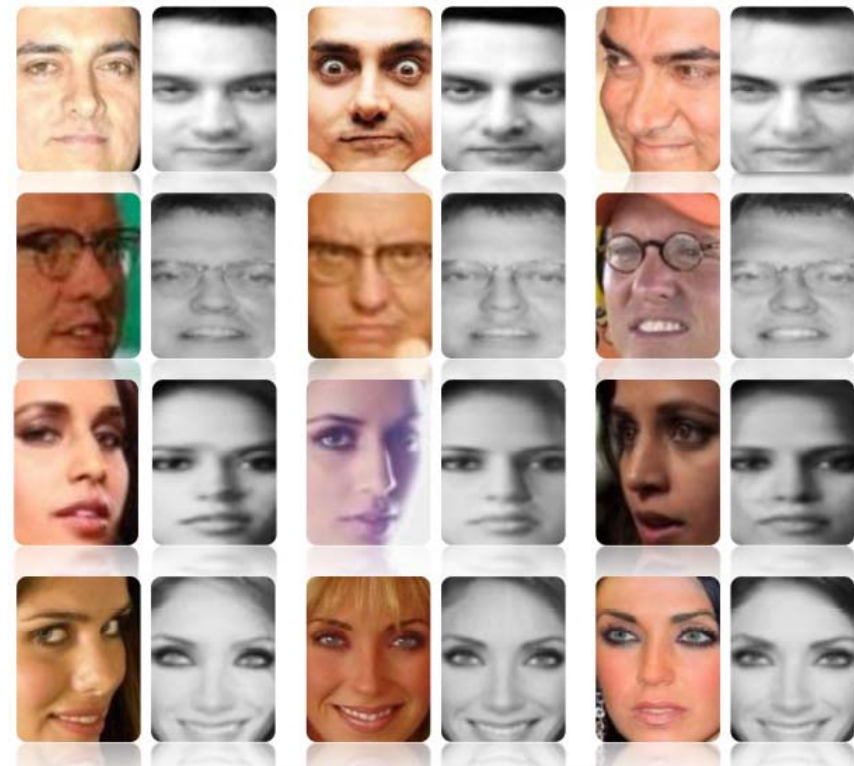
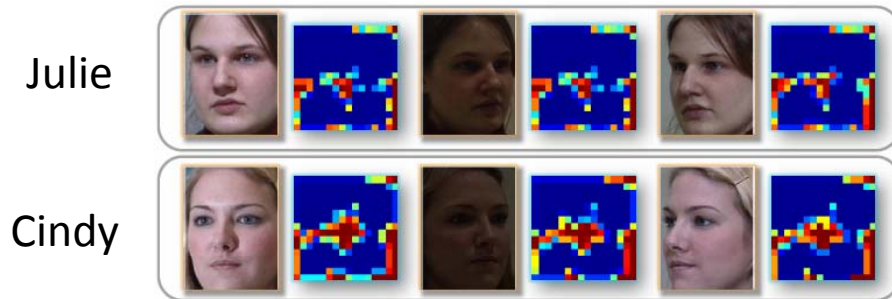
[2] Cao, Wipf, Wen, Duan, and Sun. A practical transfer learning algorithm for face verification. ICCV, 2013.

[3] Taigman, Yang, Ranzato, and Wolf. DeepFace: Closing the gap to human-level performance in face verification. CVPR, 2014.

[4] Sun, Wang, and Tang. Deep learning face representation from predicting 10,000 classes. CVPR, 2014.

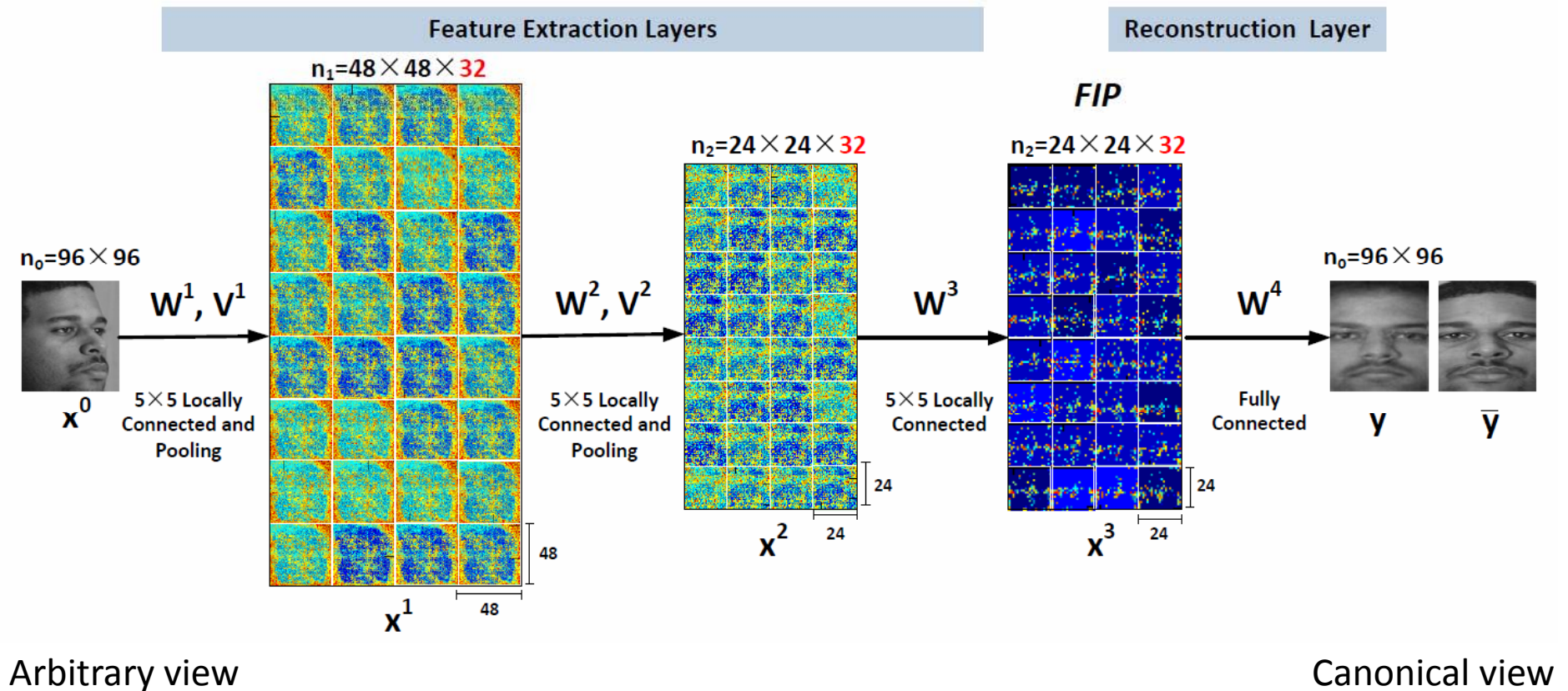
[5] Y. Sun, Y. Chen, X. Wang, and X. Tang. Deep Learning Face Representation by Joint Identification-Verification. NIPS, 2014.

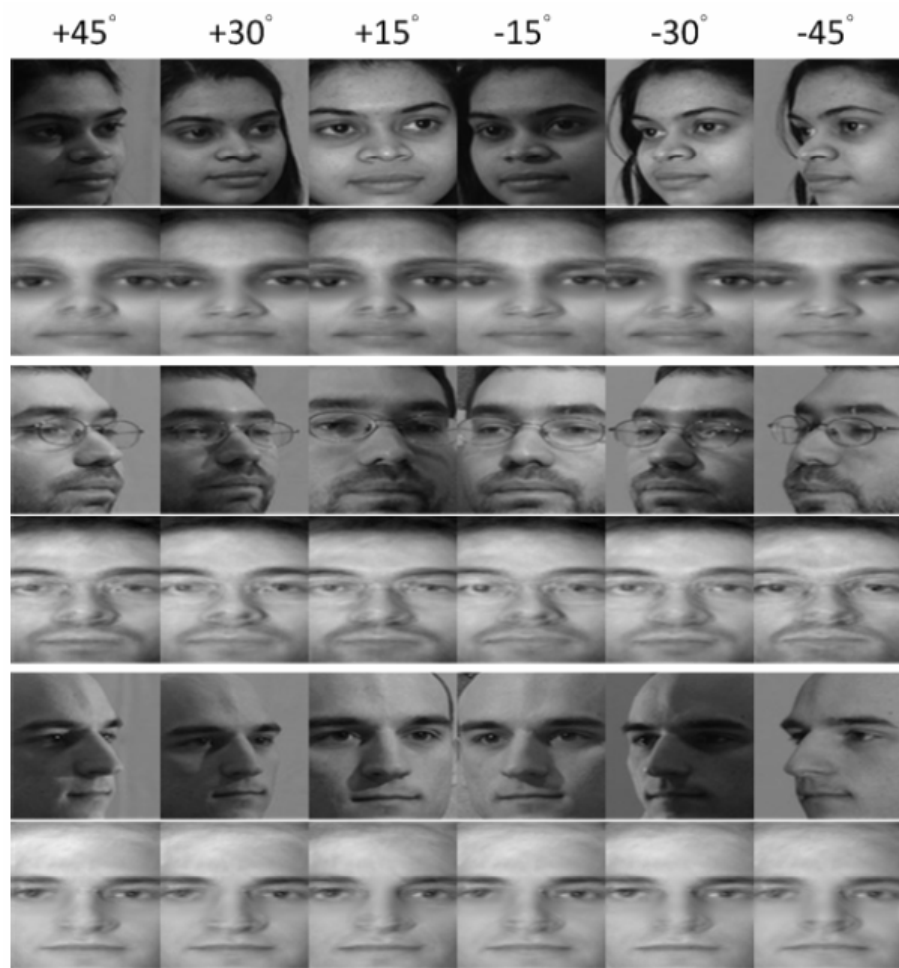
Learning face representation from recovering canonical-view face images



Reconstruction examples from LFW

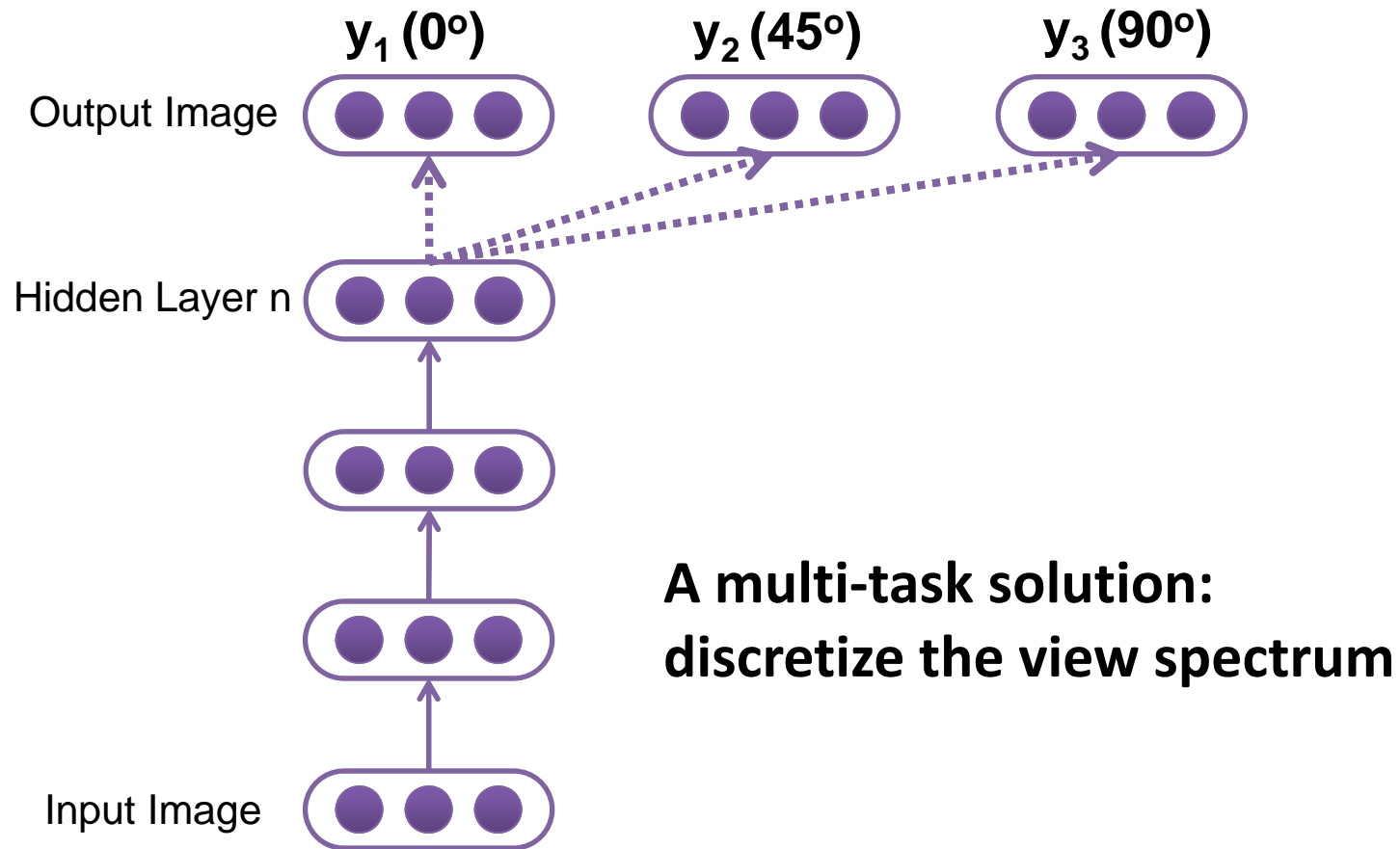
- Disentangle factors through feature extraction over multiple layers
- No 3D model; no prior information on pose and lighting condition
- Model multiple complex transforms
- Reconstructing the whole face is a much stronger supervision than predicting 0/1 class label





It is still not a 3D representation yet

Can we reconstruct all the views?



1. The number of views to be reconstructed is predefined, equivalent to the number of tasks
2. Cannot reconstruct views not presented in the training set
3. Encounters problems when the training data of different views are unbalanced
4. Model complexity increases as the number of views

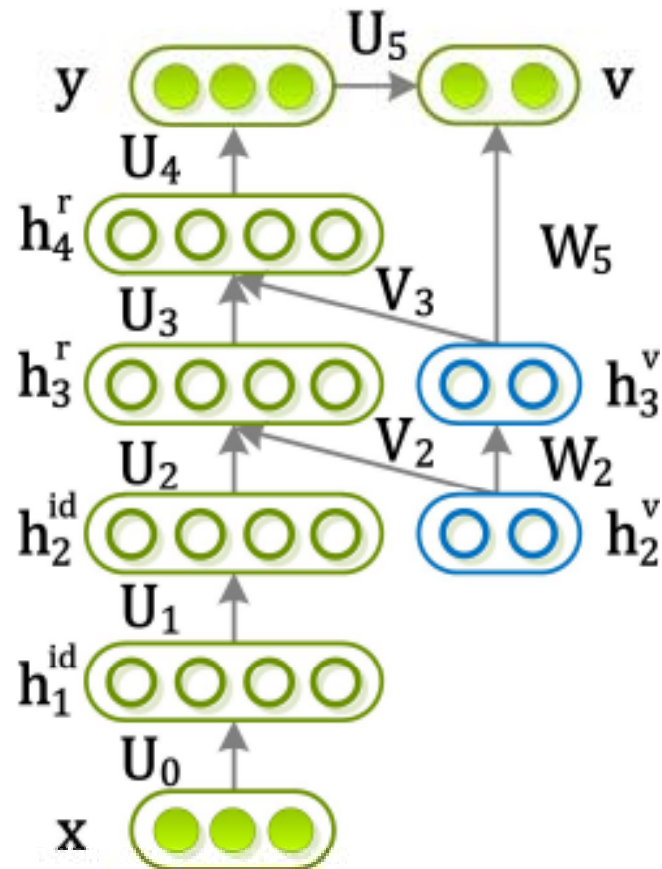
Deep learning multi-view representation from 2D images

- Given an image under arbitrary view, its viewpoint can be estimated and its full spectrum of views can be reconstructed
- Continuous view representation
- Identity and view represented by different sets of neurons



Z. Zhu, P. Luo, X. Wang, and X. Tang, "Deep Learning and Disentangling Face Representation by Multi-View Perception," NIPS 2014.

Network is composed of deterministic neurons and random neurons



x and y are input and output images of the same identity but in different views;

v is the view label of the output image;

h^{id} are neurons encoding identity features

h^v are neurons encoding view features

h^r are neurons encoding features to reconstruct the output images

Deep Learning by EM

- EM updates on the probabilistic model are converted to forward and backward propagation

$$\mathcal{L}(\Theta, \Theta^{old}) = \sum_{\mathbf{h}^v} p(\mathbf{h}^v | \mathbf{y}, \mathbf{v}; \Theta^{old}) \log p(\mathbf{y}, \mathbf{v}, \mathbf{h}^v | \mathbf{h}^{id}; \Theta)$$

- E-step: proposes s samples of \mathbf{h}

$$\mathbf{h}_s^v \sim \mathcal{U}(0, 1)$$

$$w_s = p(\mathbf{y}, \mathbf{v} | \mathbf{h}_s^v; \Theta^{old})$$

- M-step: compute gradient refer to \mathbf{h} with largest w_s

$$\frac{\partial \mathcal{L}(\Theta)}{\partial \Theta} \simeq \frac{\partial}{\partial \Theta} \left\{ w_s \left(\log p(\mathbf{v} | \mathbf{y}, \mathbf{h}_s^v) + \log p(\mathbf{y} | \mathbf{h}^{id}, \mathbf{h}_s^v) \right) \right\}$$

	Avg.	0°	-15°	+15°	-30°	+30°	-45°	+45°	-60°	+60°
Raw Pixels+LDA	36.7	81.3	59.2	58.3	35.5	37.3	21.0	19.7	12.8	7.63
LBP [1]+LDA	50.2	89.1	77.4	79.1	56.8	55.9	35.2	29.7	16.2	14.6
Landmark LBP [6]+LDA	63.2	94.9	83.9	82.9	71.4	68.2	52.8	48.3	35.5	32.1
CNN+LDA	58.1	64.6	66.2	62.8	60.7	63.6	56.4	57.9	46.4	44.2
FIP [28]+LDA	72.9	94.3	91.4	90.0	78.9	82.5	66.1	62.0	49.3	42.5
RL [28]+LDA	70.8	94.3	90.5	89.8	77.5	80.0	63.6	59.5	44.6	38.9
MTL+RL+LDA	74.8	93.8	91.7	89.6	80.1	83.3	70.4	63.8	51.5	50.2
MVP _{h₁} ^{id} +LDA	61.5	92.5	85.4	84.9	64.3	67.0	51.6	45.4	35.1	28.3
MVP _{h₂} ^{id} +LDA	79.3	95.7	93.3	92.2	83.4	83.9	75.2	70.6	60.2	60.0
MVP _{h₃} ^r +LDA	72.6	91.0	86.7	84.1	74.6	74.2	68.5	63.8	55.7	56.0
MVP _{h₄} ^r +LDA	62.3	83.4	77.3	73.1	62.0	63.9	57.3	53.2	44.4	46.9

Face recognition accuracies across views and illuminations on the Multi-PIE dataset. The first and the second best performances are in bold.

- [1] T. Ahonen, A. Hadid, and M. Pietikainen. Face description with local binary patterns: Application to face recognition. *TPAMI*, 28:2037–2041, 2006.
- [6] Dong Chen, Xudong Cao, Fang Wen, and Jian Sun. Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification. In *CVPR*, 2013.
- [28] Z. Zhu, P. Luo, X. Wang, and X. Tang. Deep learning identity preserving face space. In *ICCV*, 2013.

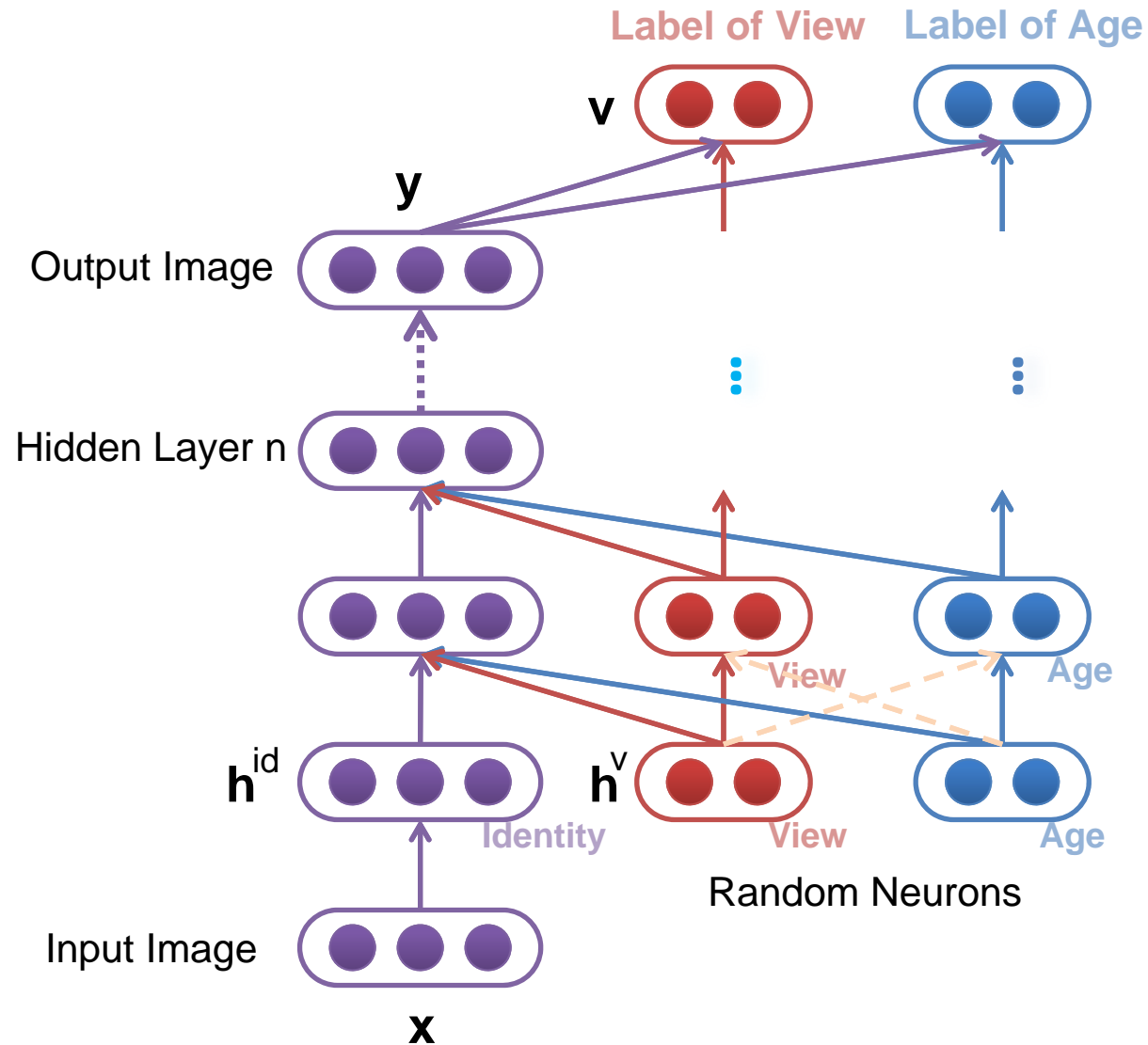
Deep Learning Multi-view Representation from 2D Images

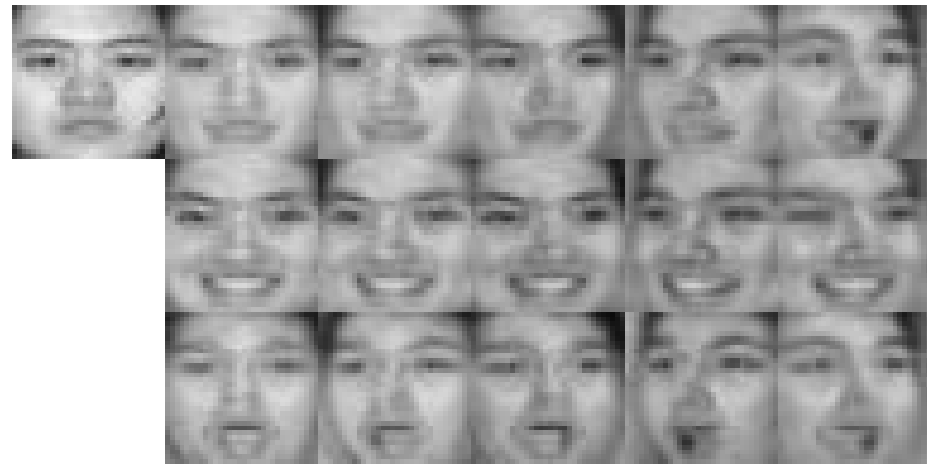
- Interpolate and predict images under viewpoints unobserved in the training set



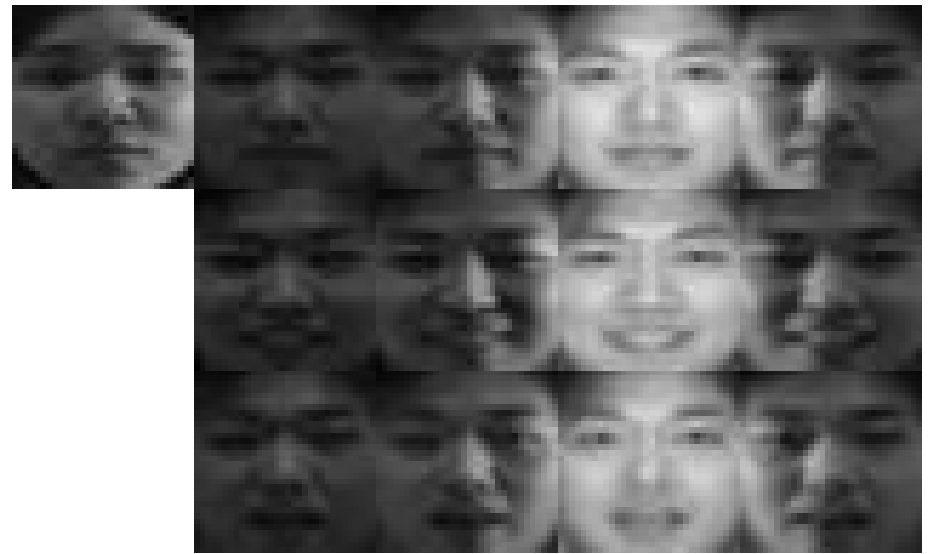
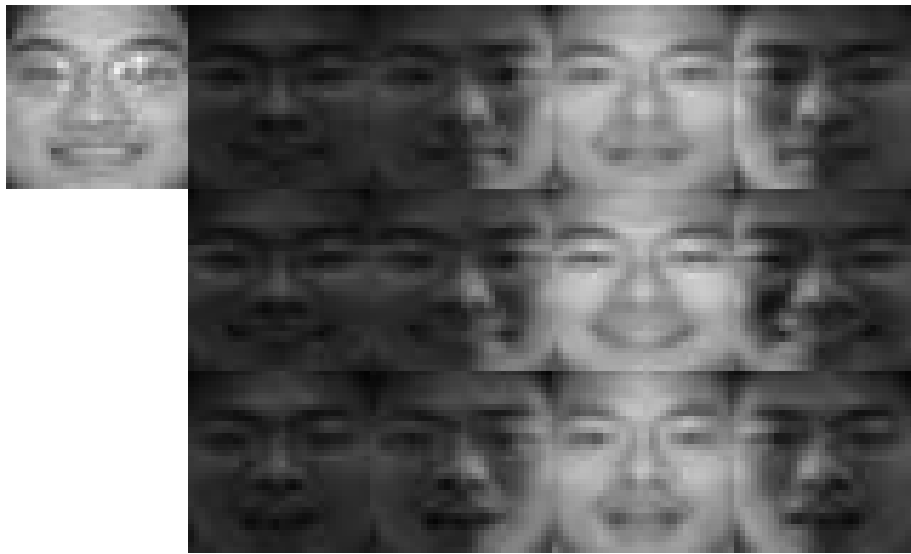
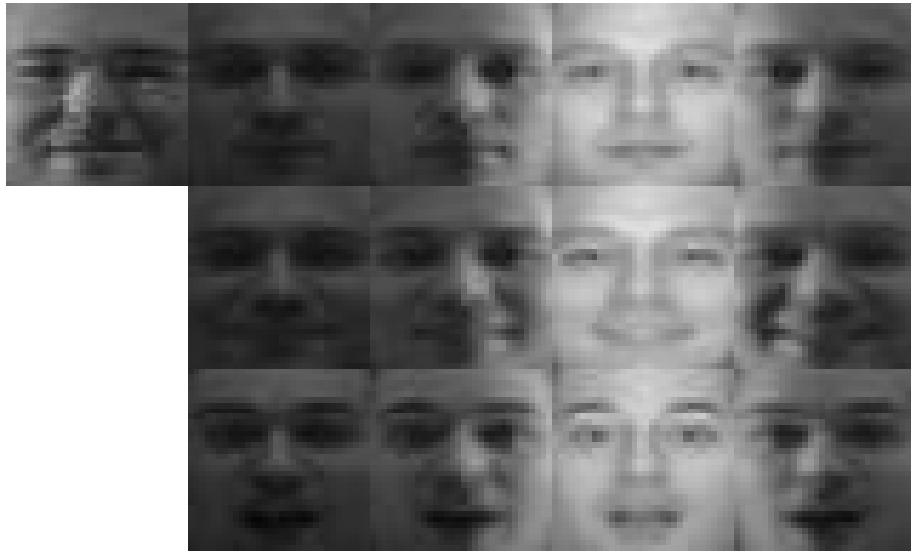
The training set only has viewpoints of 0° , 30° , and 60° . (a): the reconstructed images under 15° and 45° when the input is taken under 0° . (b) The input images are under 15° and 45° .

Generalize to other facial factors





Face reconstruction across poses and expressions



Face reconstruction across lightings and expressions

Learn face representations from

face verification, identification, multi-view reconstruction

Properties of face representations

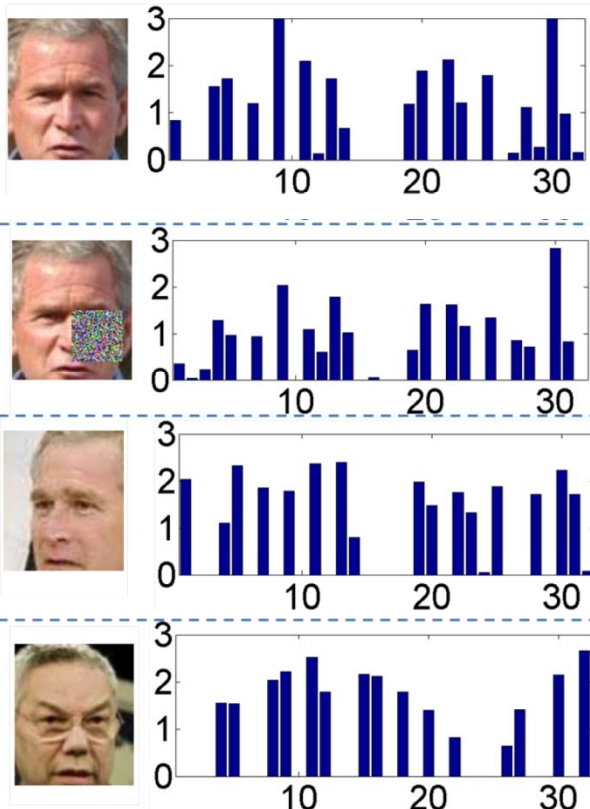
sparseness, selectiveness, robustness

Applications of face representations

face localization, attribute recognition

Y. Sun, X. Wang, and X. Tang, “Deeply learned face representations are sparse, selective, and robust,” CVPR 2015

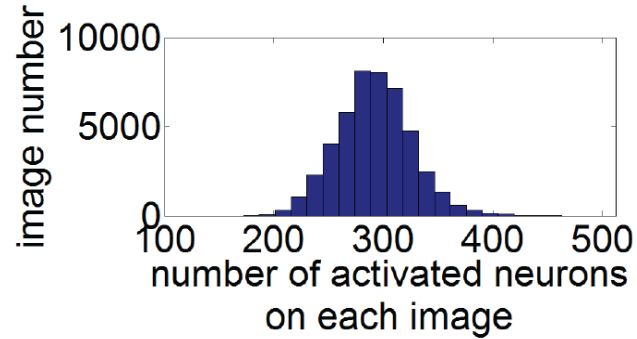
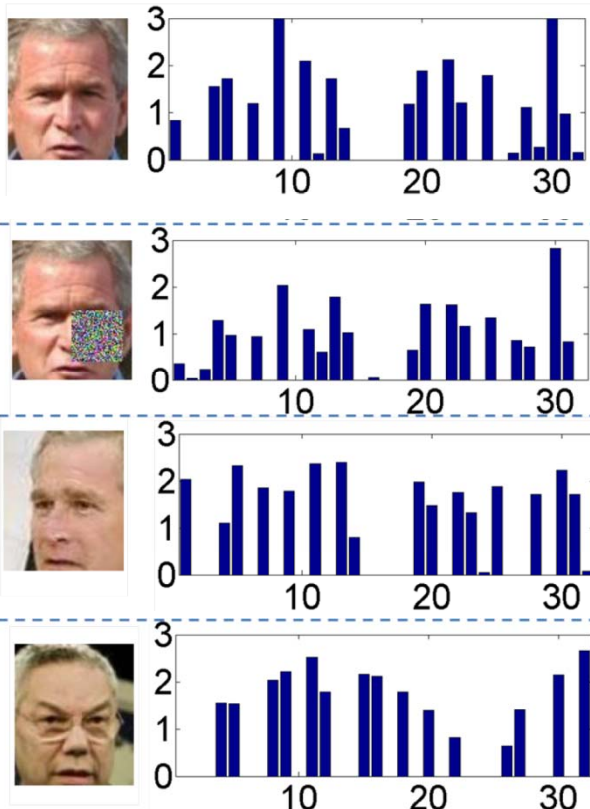
Deeply learned features are moderately sparse



- The **binary codes** on activation patterns are very effective on face recognition
- Save storage and speedup face search dramatically
- Activation patterns are more important than activation magnitudes in face recognition

	Joint Bayesian (%)	Hamming distance (%)
Combined model (real values)	99.47	n/a
Combined model (binary code)	99.12	97.47

Deeply learned features are moderately sparse



1	0	1	1	0	0
0	1	0	0	1	1

6

Moderately sparse

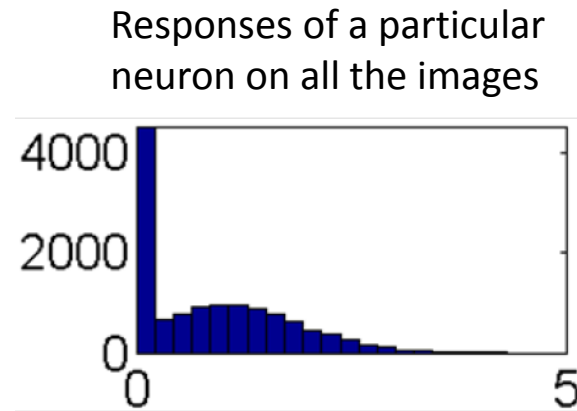
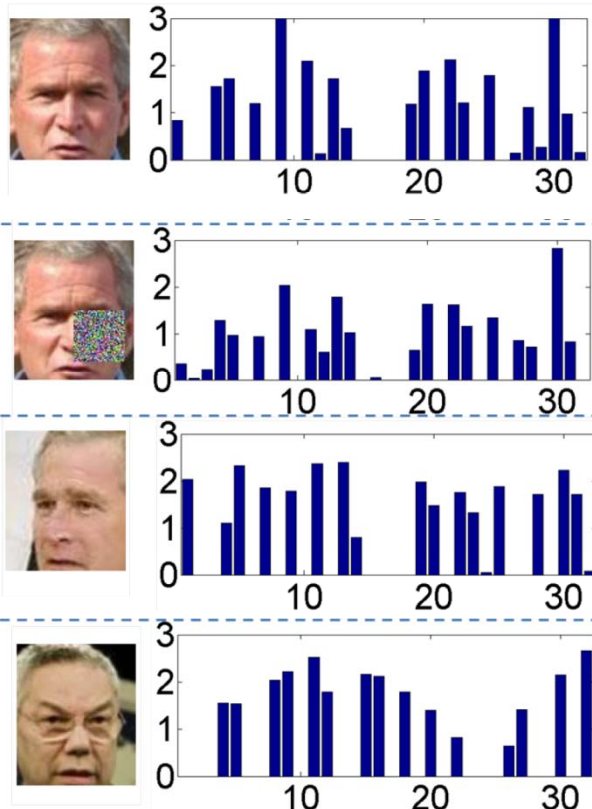
1	0	0	0	0	0
0	1	0	0	0	0

2

Highly sparse

- For an input image, about half of the neurons are activated
 - ✓ Maximize the Hamming distance between images

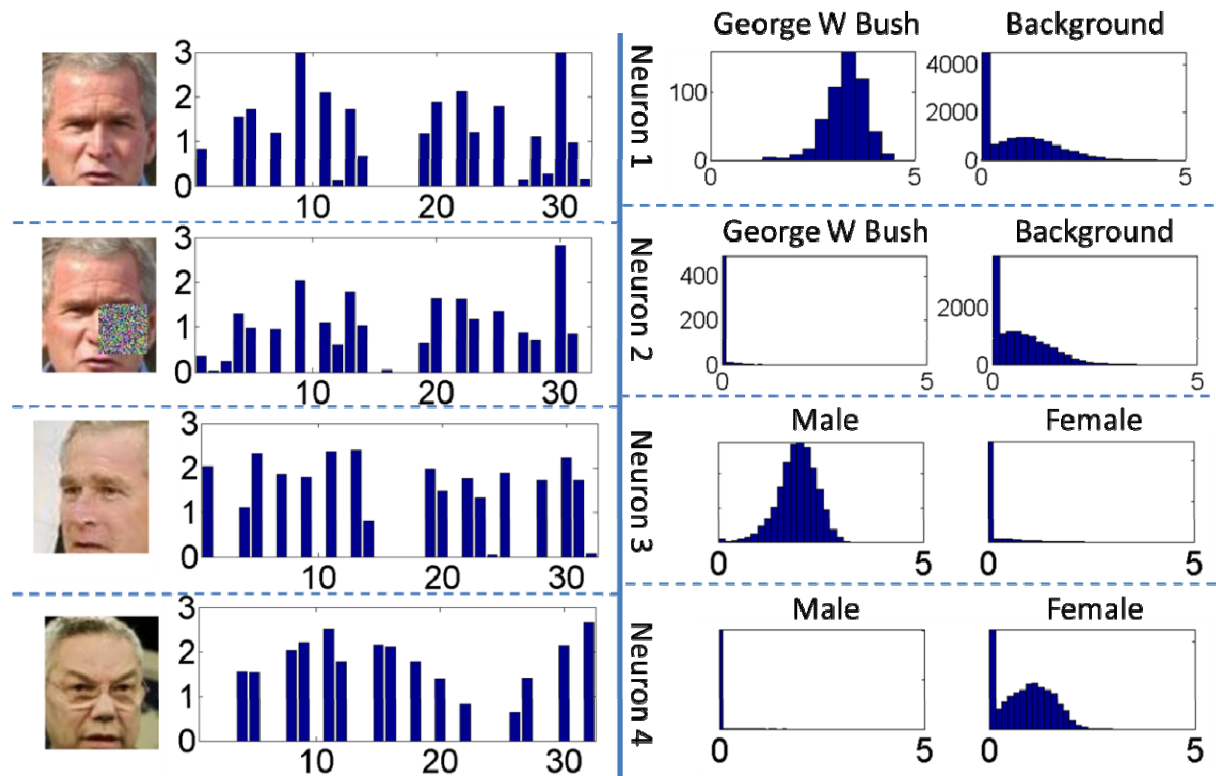
Deeply learned features are moderately sparse



- An neuron has response on about half of the images
 - ✓ Maximize the discriminative power (entropy) of a neuron on describing the image set

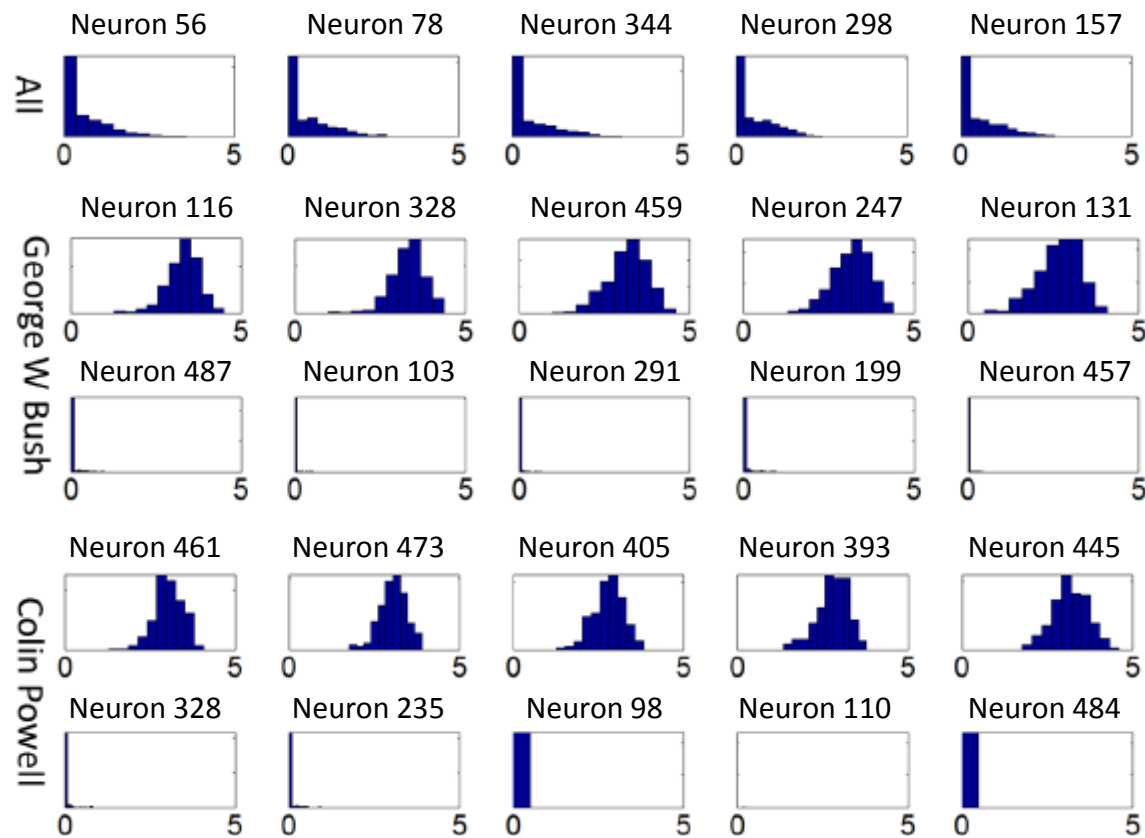
Deeply learned features are selective to identities and attributes

- With a single neuron, DeepID2 reaches 97% recognition accuracy for some identity and attribute

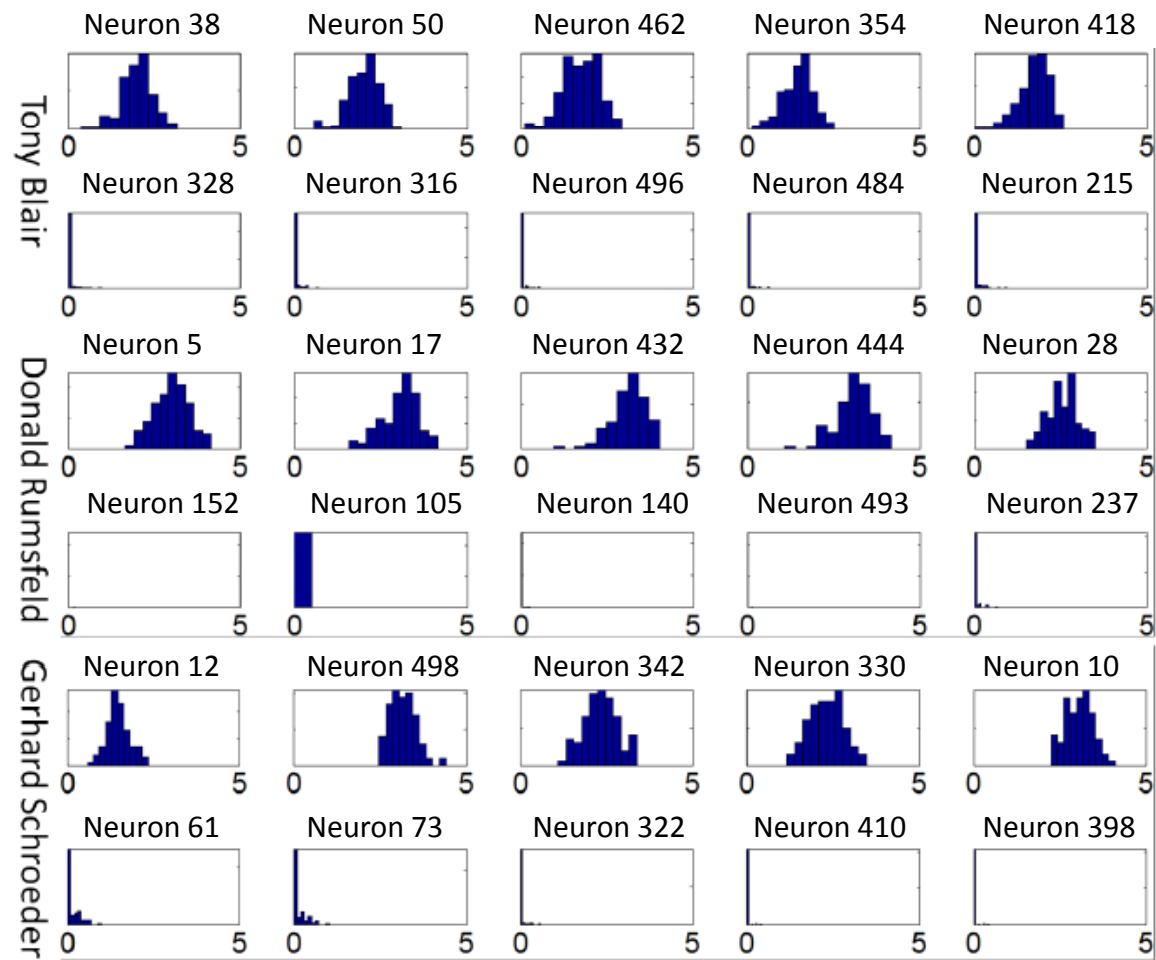


Deeply learned features are selective to identities and attributes

- Excitatory and inhibitory neurons (on identities)

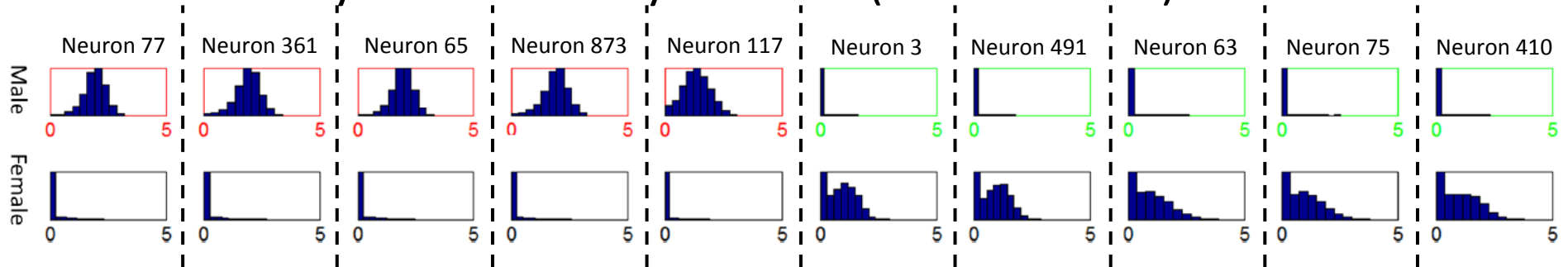


Histograms of neural activations over identities with the most images in LFW

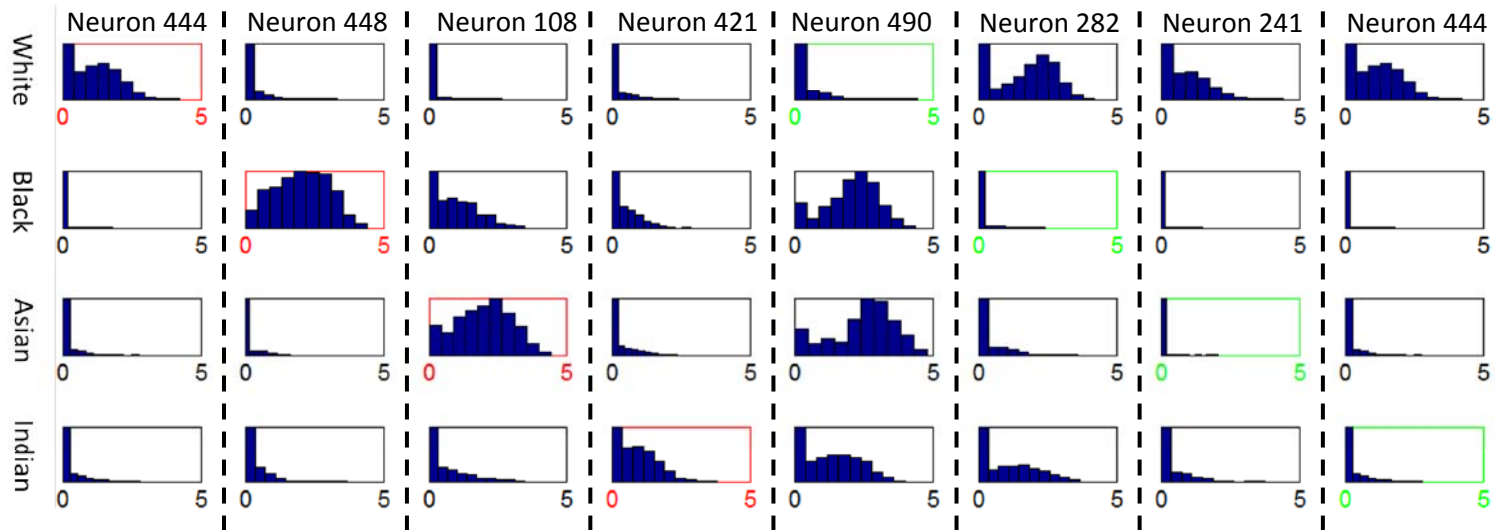


Deeply learned features are selective to identities and attributes

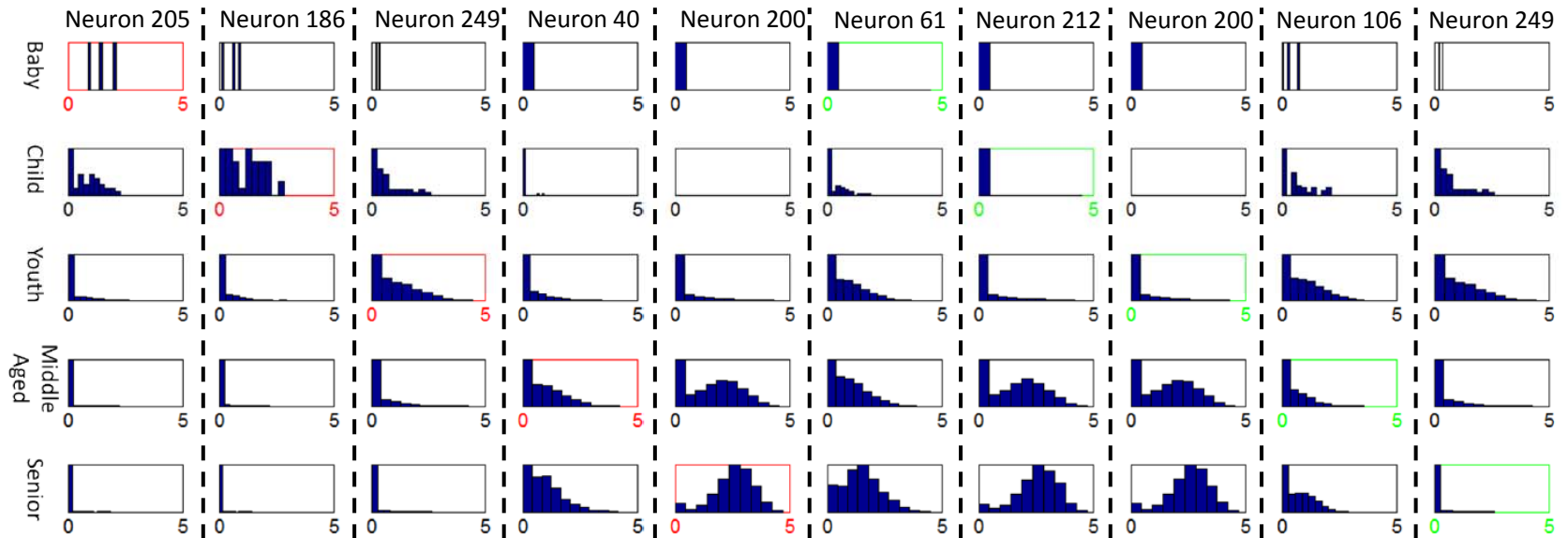
- Excitatory and inhibitory neurons (on attributes)



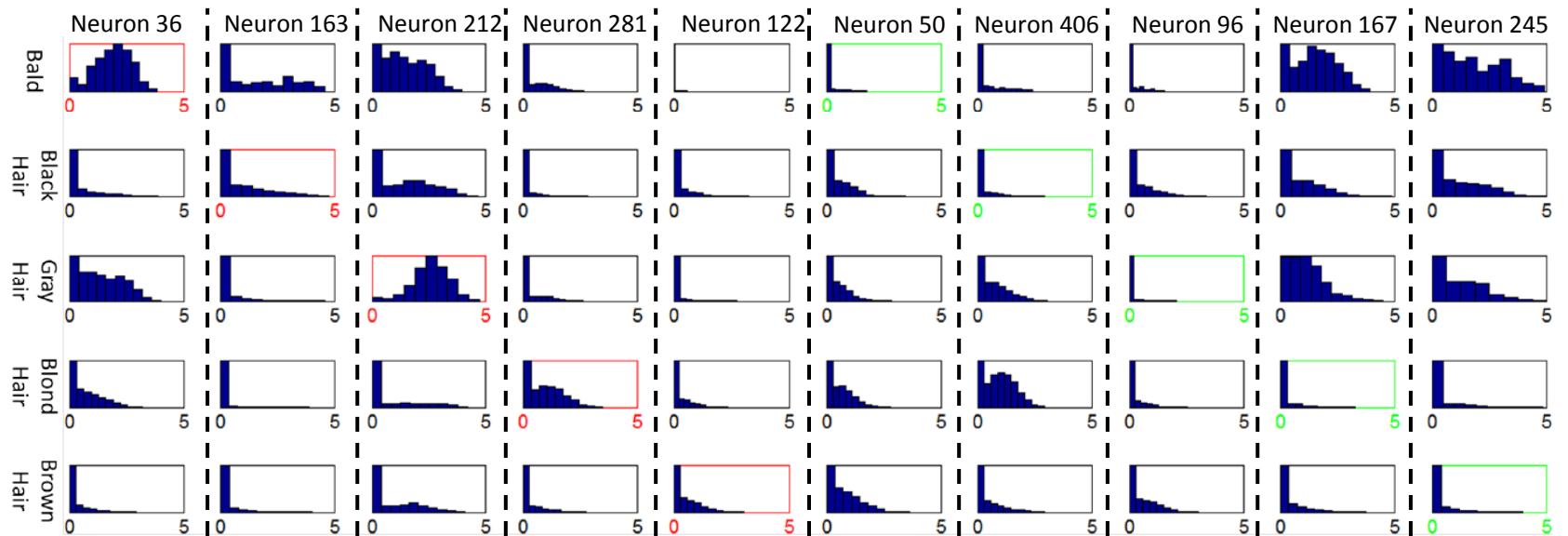
Histograms of neural activations over gender-related attributes (Male and Female)



Histograms of neural activations over race-related attributes (White, Black, Asian and India)



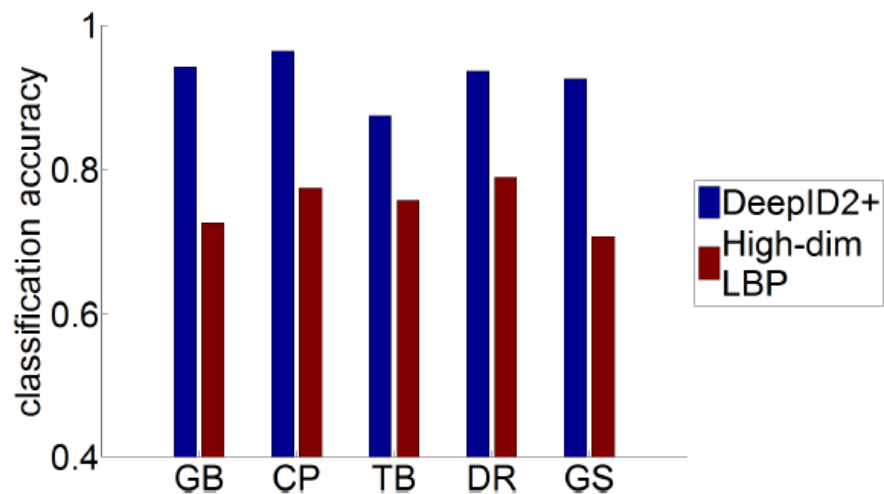
Histogram of neural activations over age-related attributes (Baby, Child, Youth, Middle Aged, and Senior)



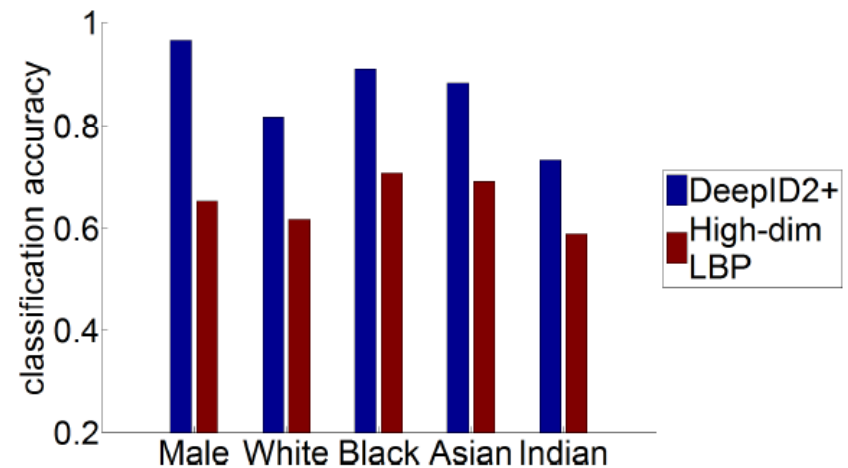
Histogram of neural activations over hair-related attributes (Bald, Black Hair, Gray Hair, Blond Hair, and Brown Hair)

Deeply learned features are selective to identities and attributes

- With a single neuron, DeepID2 reaches 97% recognition accuracy for some identity and attribute



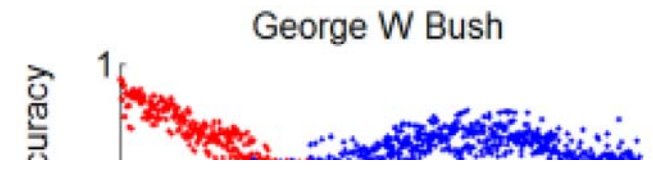
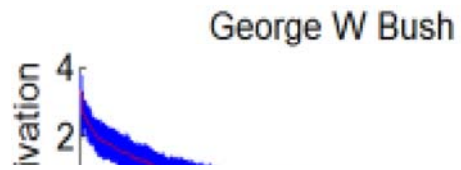
Identity classification accuracy on LFW with one single DeepID2+ or LBP feature. GB, CP, TB, DR, and GS are five celebrities with the most images in LFW.



Attribute classification accuracy on LFW with one single DeepID2+ or LBP feature.

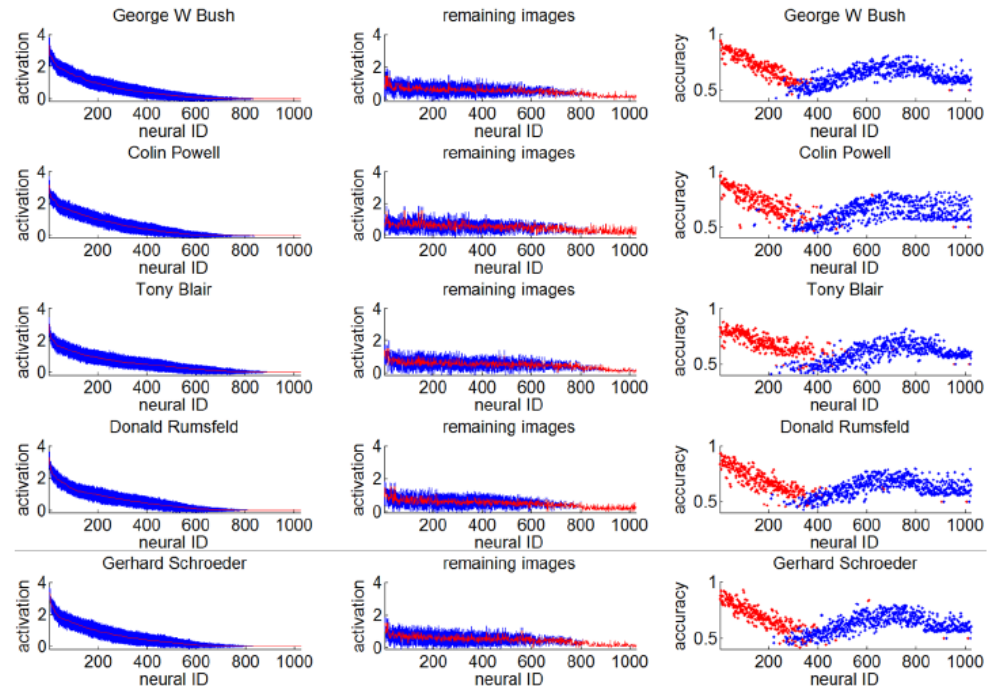
Excitatory and Inhibitory neurons

DeepID2+

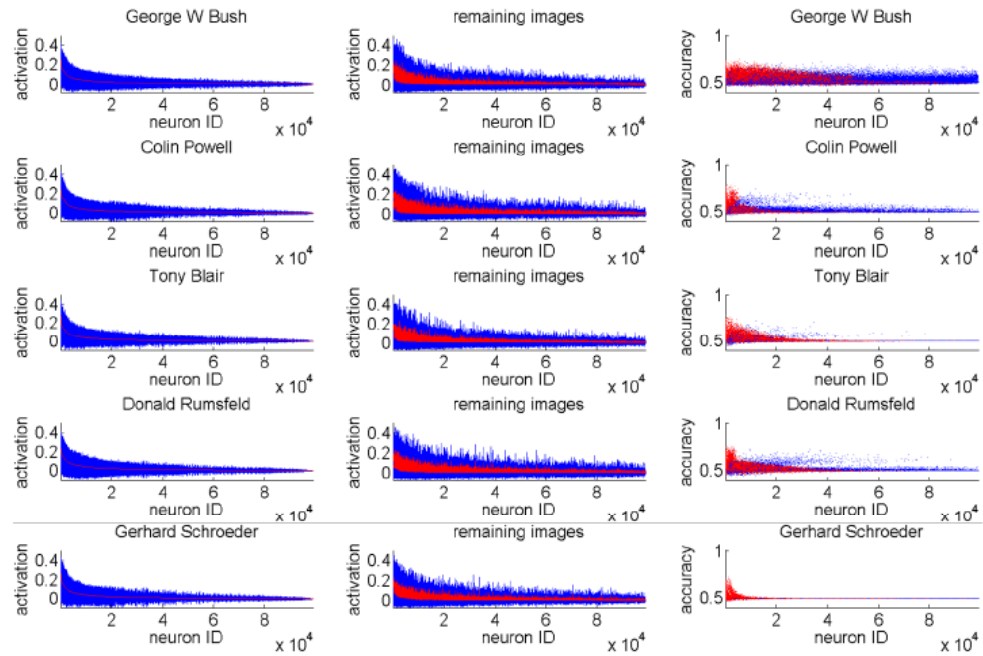


High-dim LBP

Excitatory and Inhibitory neurons

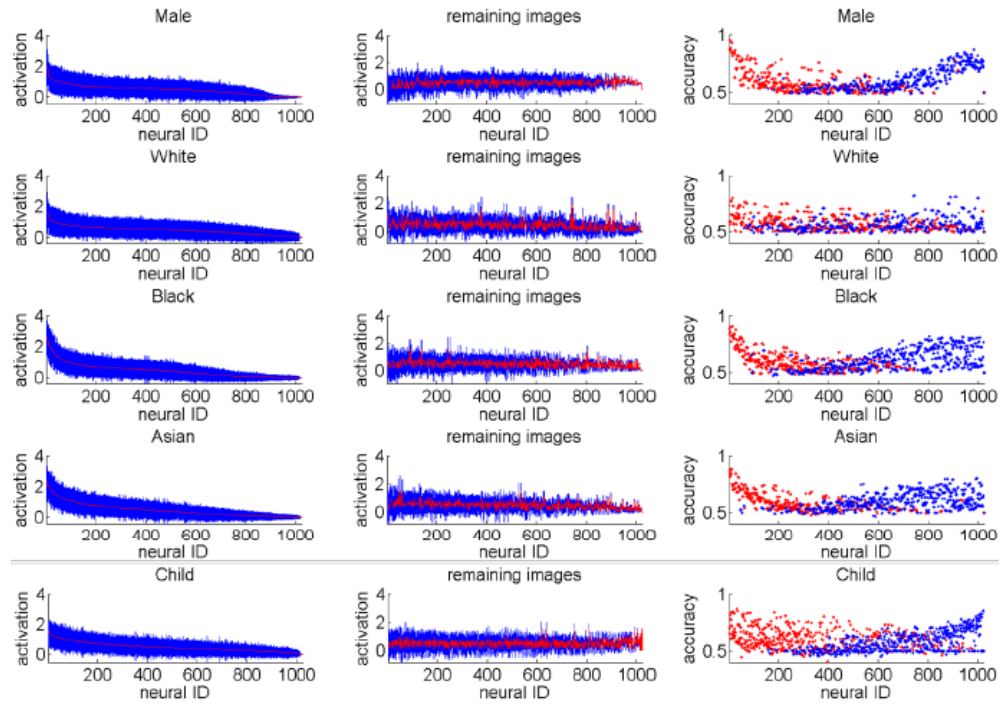


DeepID2+

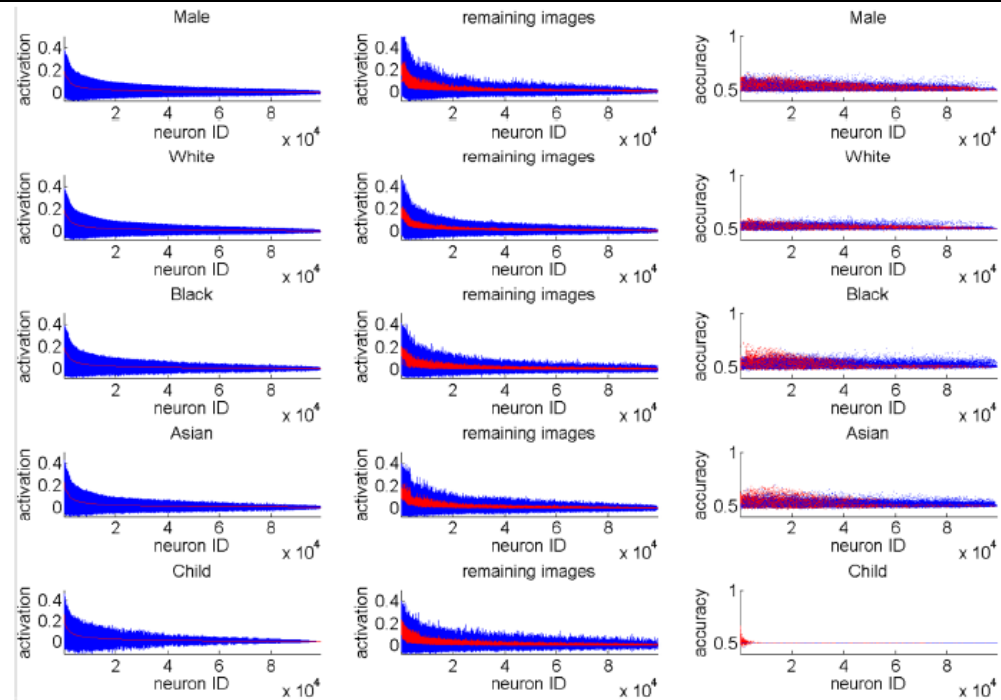


High-dim LBP

Excitatory and Inhibitory neurons



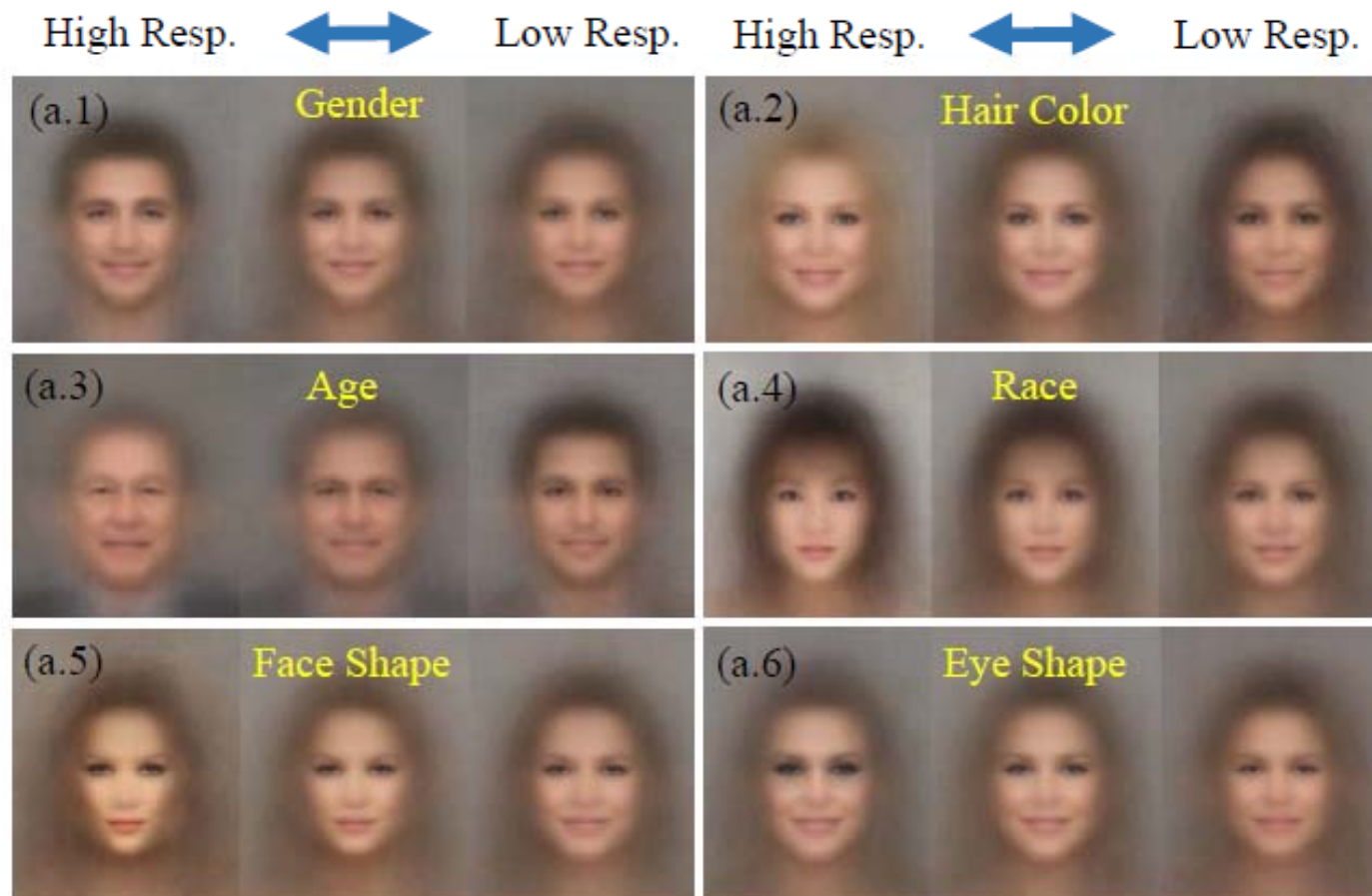
DeepID2+



High-dim LBP

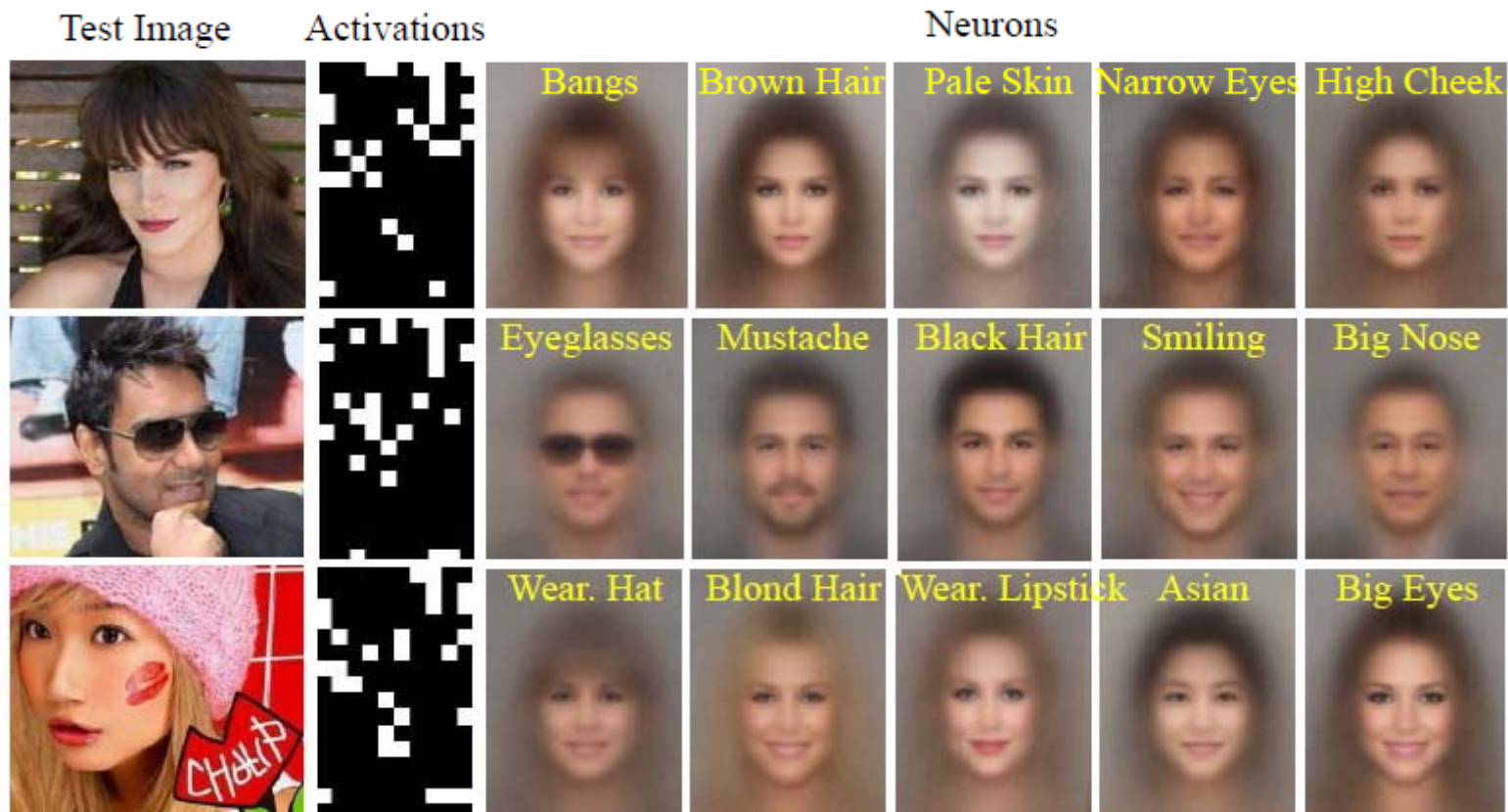
Deeply learned features are selective to identities and attributes

- Visualize the semantic meaning of each neuron



Deeply learned features are selective to identities and attributes

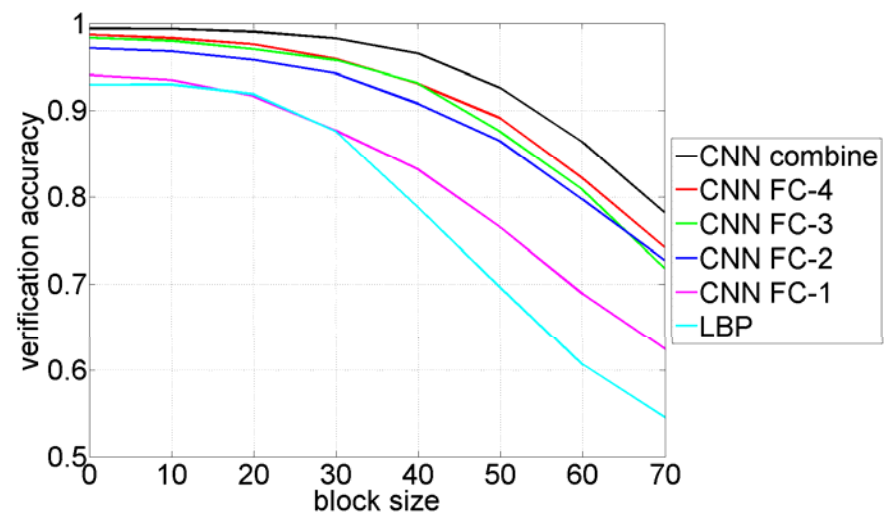
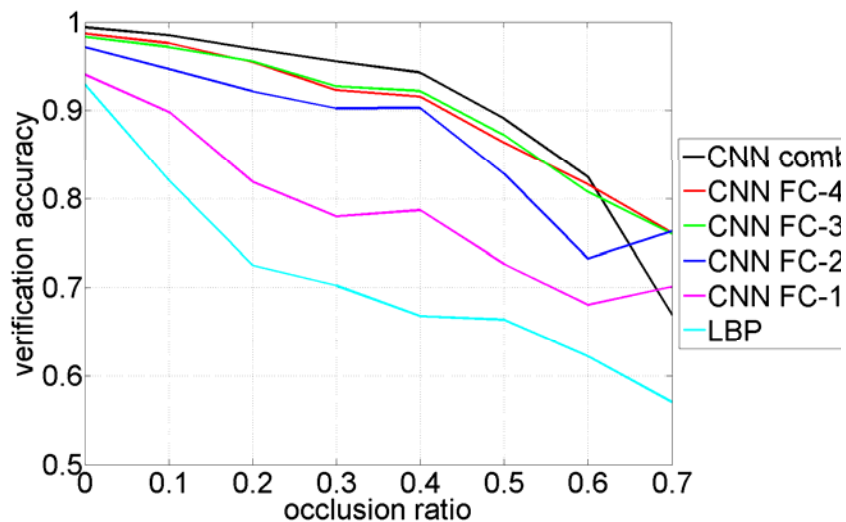
- Visualize the semantic meaning of each neuron



Neurons are ranked by their responses in descending order with respect to test images

Deeply learned features are robust to occlusions

- Global features are more robust to occlusions



Learn face representations from

face verification, identification, multi-view reconstruction

Properties of face representations

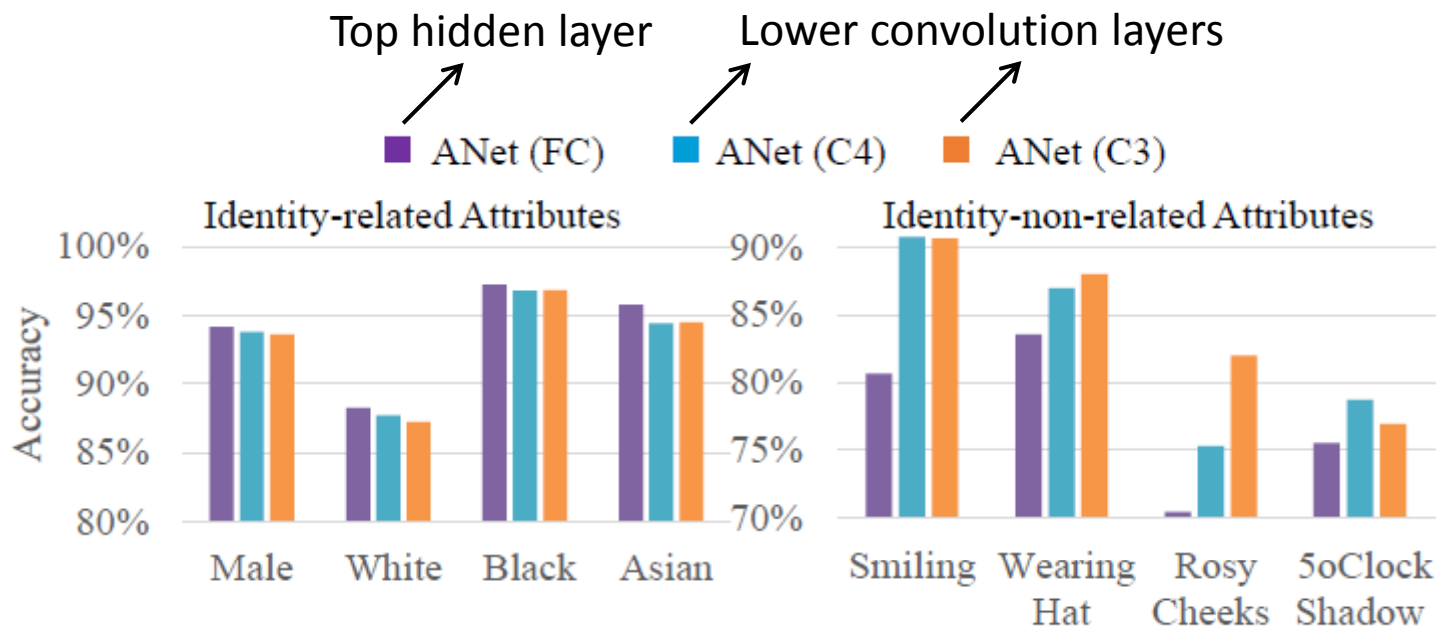
sparseness, selectiveness, robustness

Applications of face representations

face localization, attribute recognition

DeepID2 features for attribute recognition

- Features at top layers are more effective on recognizing identity related attributes
- Features at lower layers are more effective on identity-non-related attributes

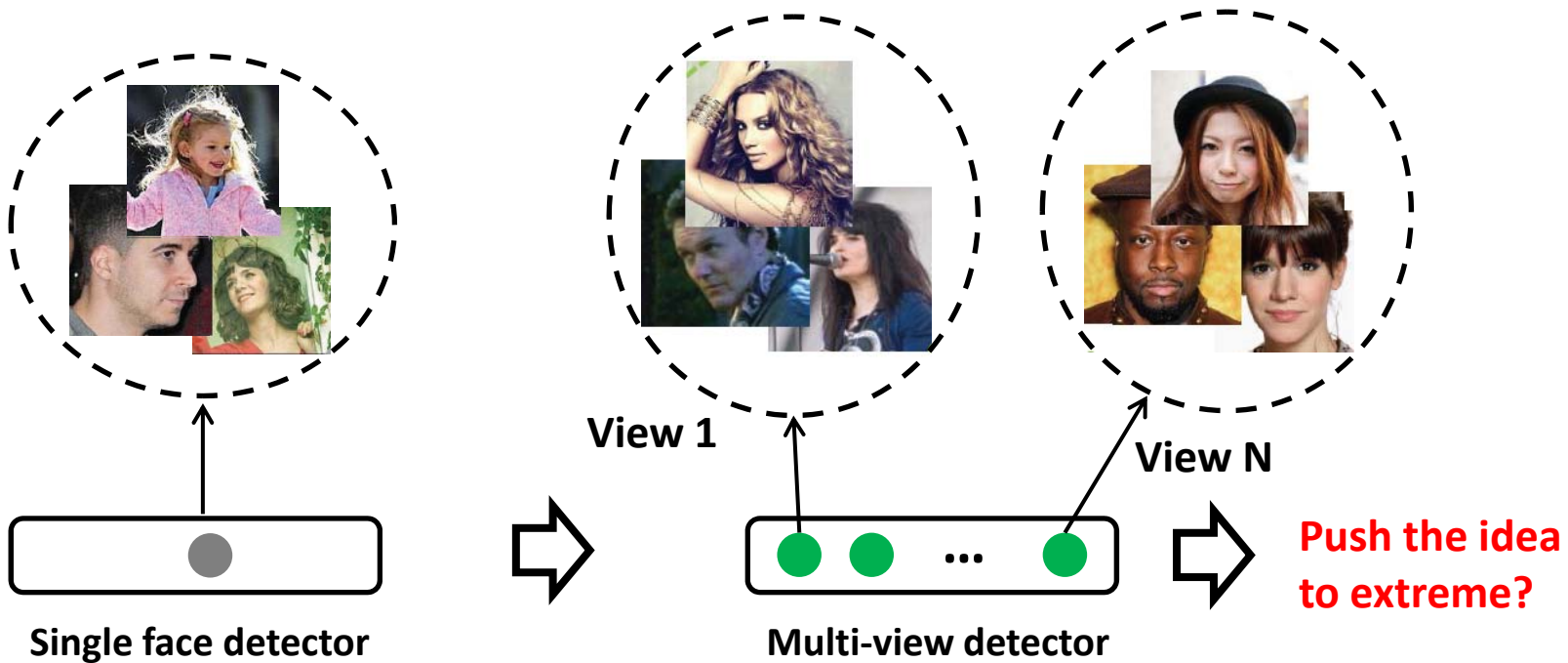


DeepID2 features for attribute recognition

- DeepID2 features can be directly used for attribute recognition
- Use DeepID2 features as initialization (pre-trained result), and then fine tune on attribute recognition
- Multi-task learning face recognition and attribute prediction does not improve performance, because face recognition is a much stronger supervision than attribute prediction
- Average accuracy on 40 attributes on CelebA and LFWA datasets

	CelebA	LFWA
FaceTracer [1] (HOG+SVM)	81	74
Training CNN from scratch with attributes	83	79
Directly use DeepID2 features	84	82
DeepID2 + fine-tuning	87	84

Features learned from face recognition can improve face localization?



Hard to handle large variety especially on views

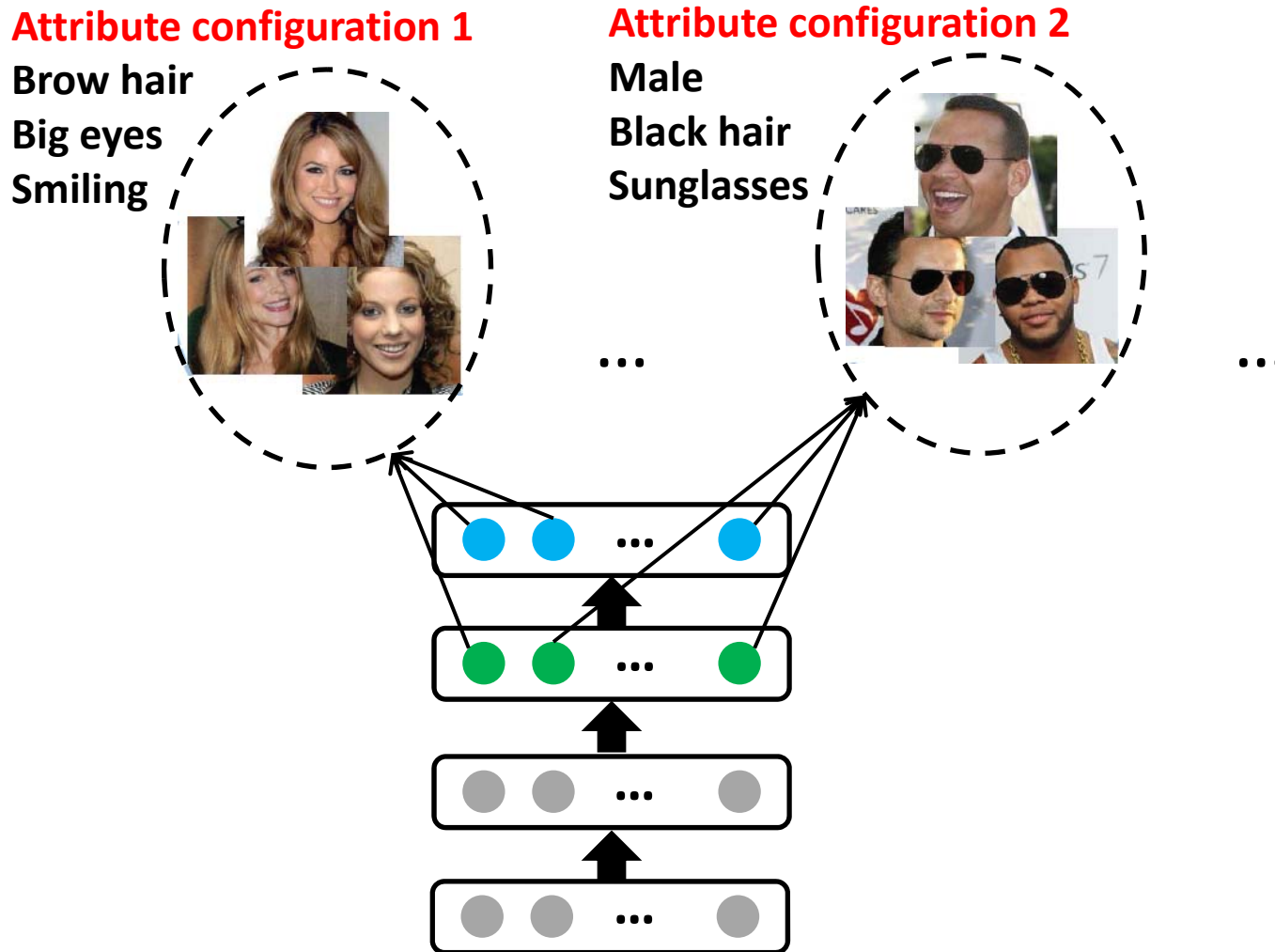
View labels are given in training; Each detector handles a view

Viewpoints → Gender, expression, race, hair style → Attributes

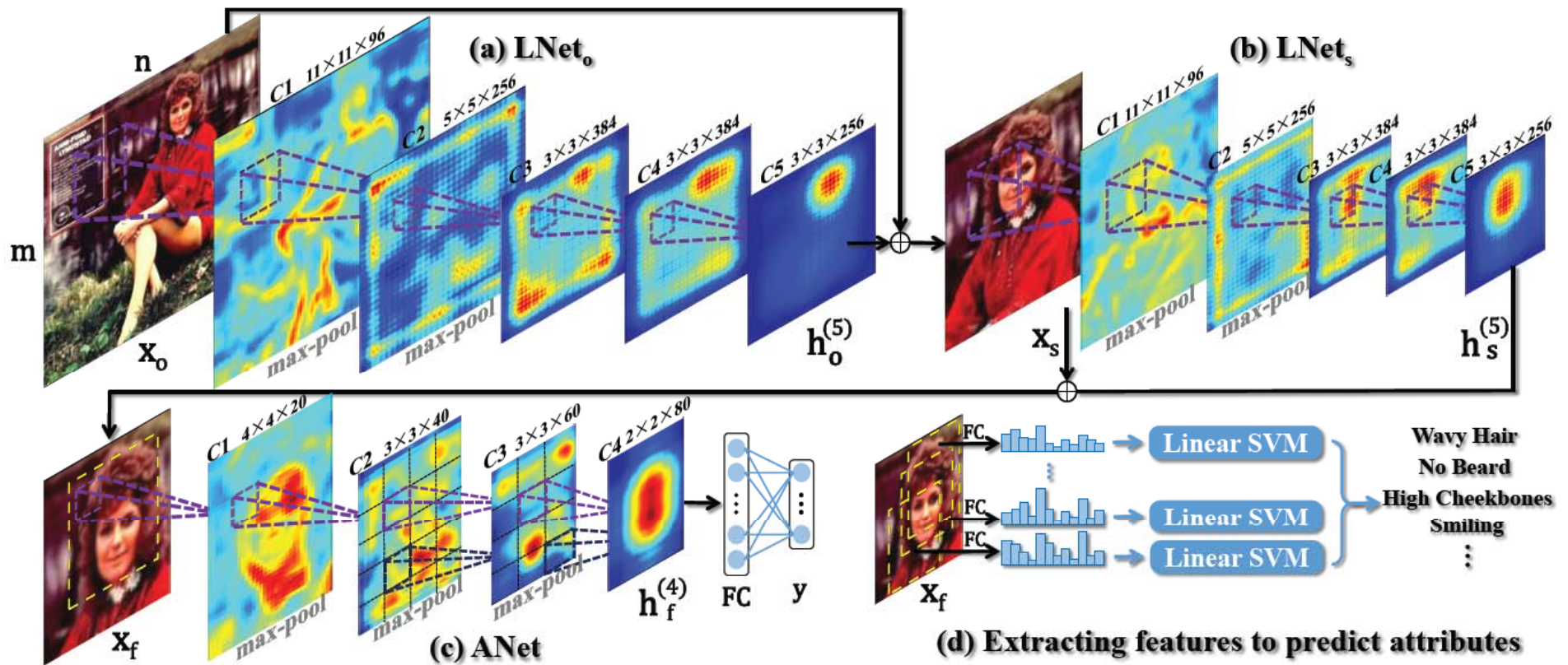
Neurons have selectiveness on attributes

A filter (or a group of filters) functions as a detector of a face attribute

When a subset of neurons are activated, they indicate existence of faces with an attribute configuration



The neurons at different layers can form many activation patterns, implying that the whole set of face images can be divided into many subsets based on attribute configurations

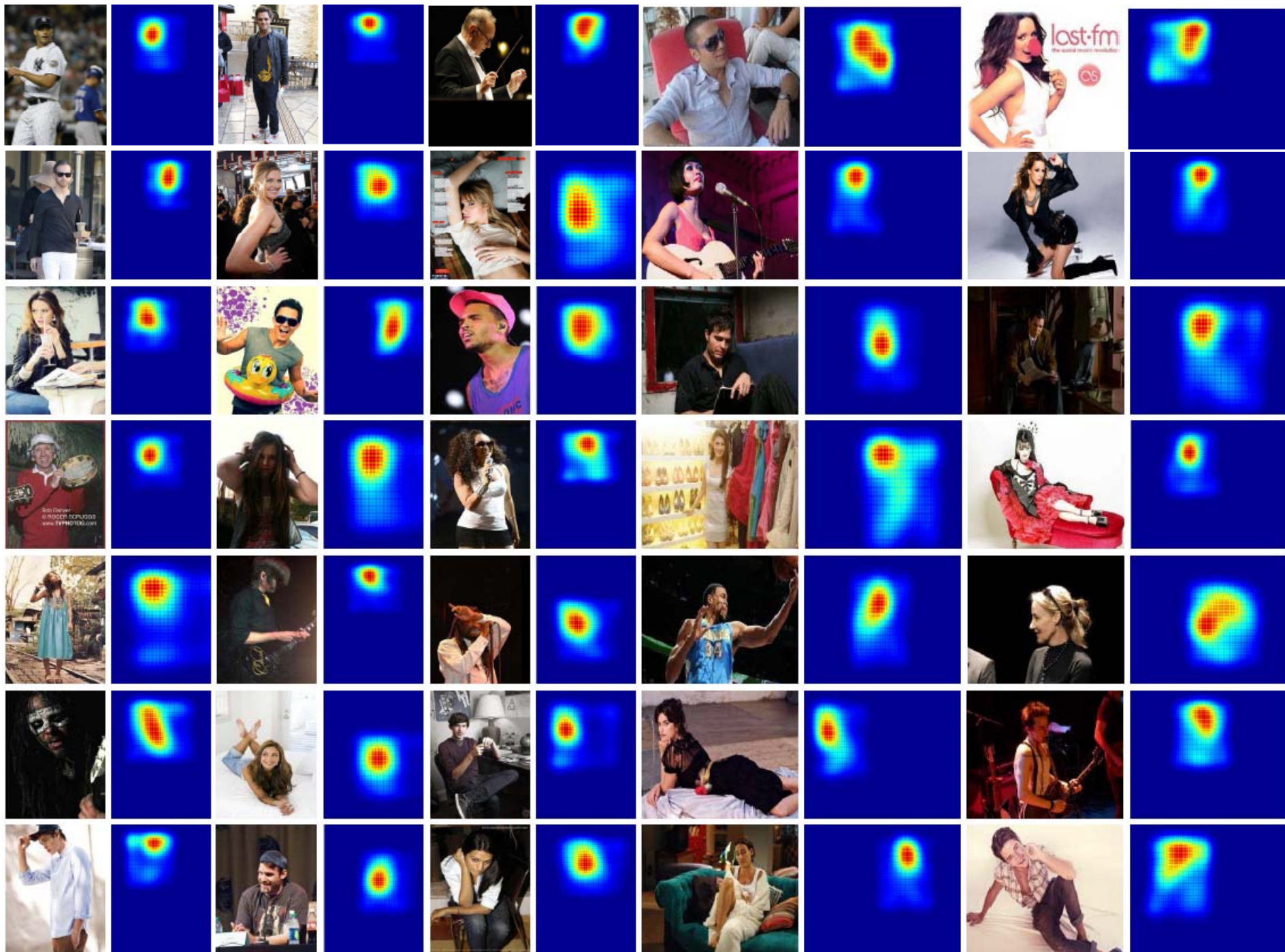


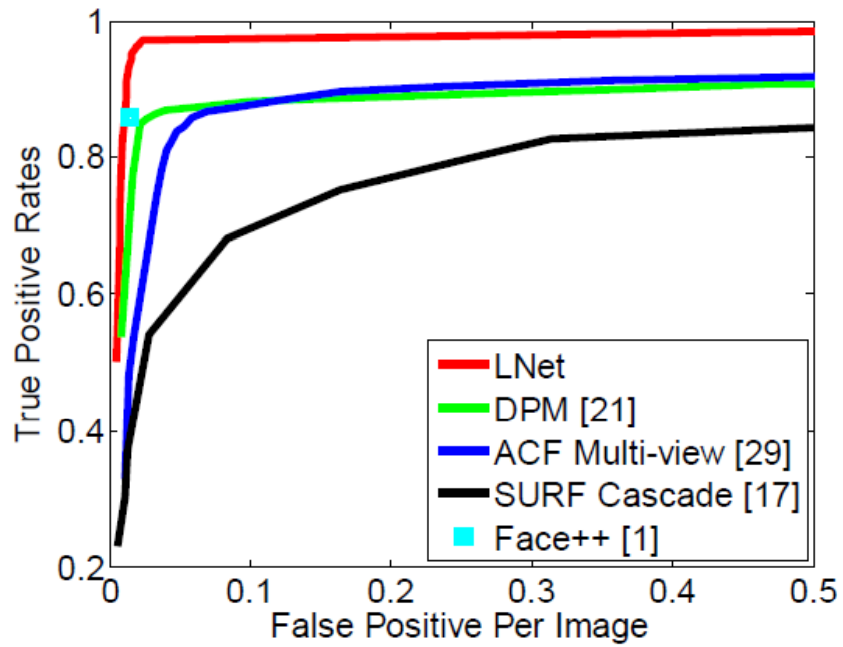
LNet localizes faces

LNet is pre-trained with face recognition and fine-tuned with attribute prediction

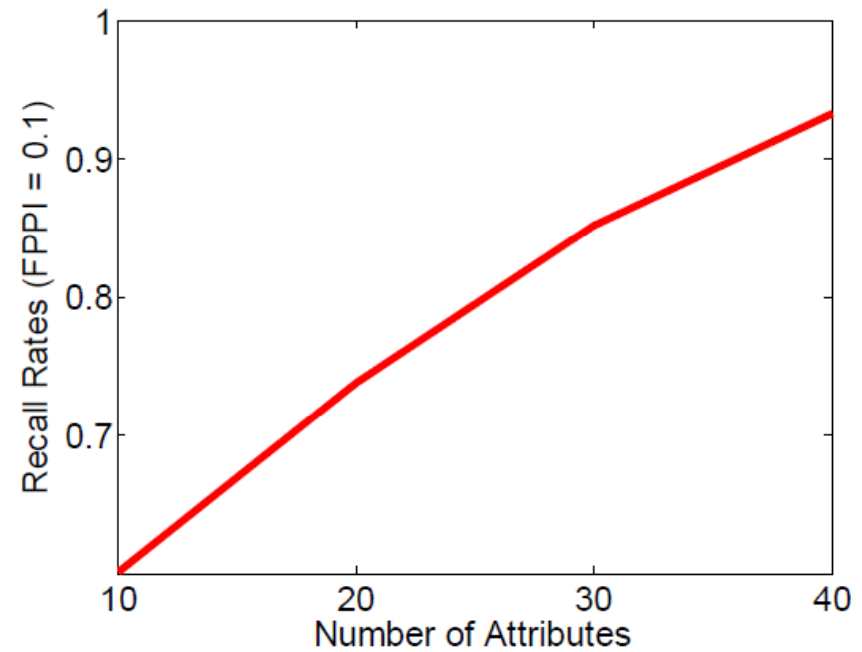
By simply averaging response maps and good face localization is achieved

Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep Learning Face Attributes in the Wild," ICCV 2015





(a)

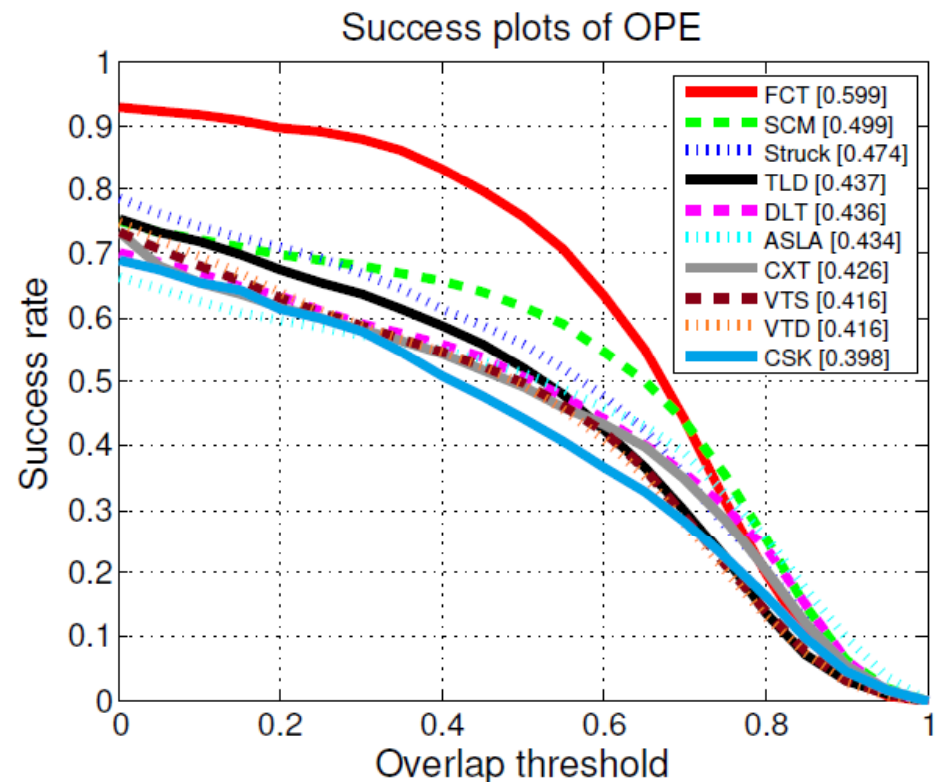
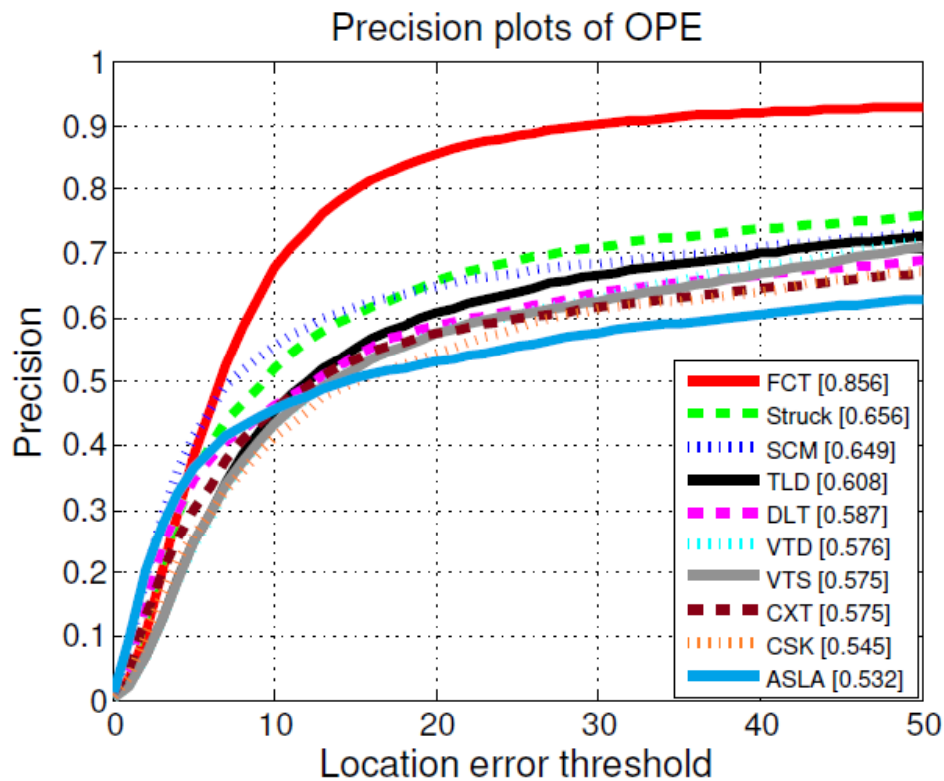


(b)

(a) ROC curves of LNet and state-of-the-art face detectors

(b) Recall rates w.r.t. number of attributes (FPPI = 0.1)

Attribute selectiveness: neurons serve as **detectors**
Identity selectiveness: neurons serve as **trackers**



L. Wang, W. Ouyang, X. Wang, and H. Lu, "Visual Tracking with Fully Convolutional Networks," ICCV 2015.

Conclusions

- Face representation can be learned from the tasks of verification, identification, and multi-view reconstruction
- Face representation can be more effectively learned from rich prediction and challenging tasks
- Deeply learned features are moderately sparse, identity and attribute selective, and robust to data corruption
- Binary neuron activation patterns are effective for face recognition than activation magnitudes
- These properties are naturally learned by DeepID2 through large-scale training
- Because of these properties, the learned face representation are effective for applications beyond face recognition, such as face localization and attribute prediction

Collaborators



Yi Sun



Ziwei Liu



Zhenyao Zhu

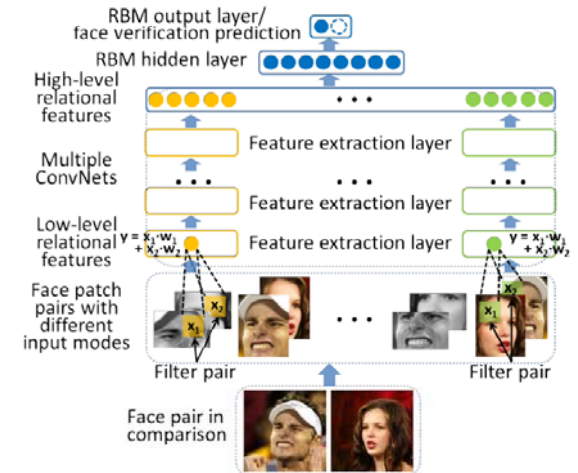
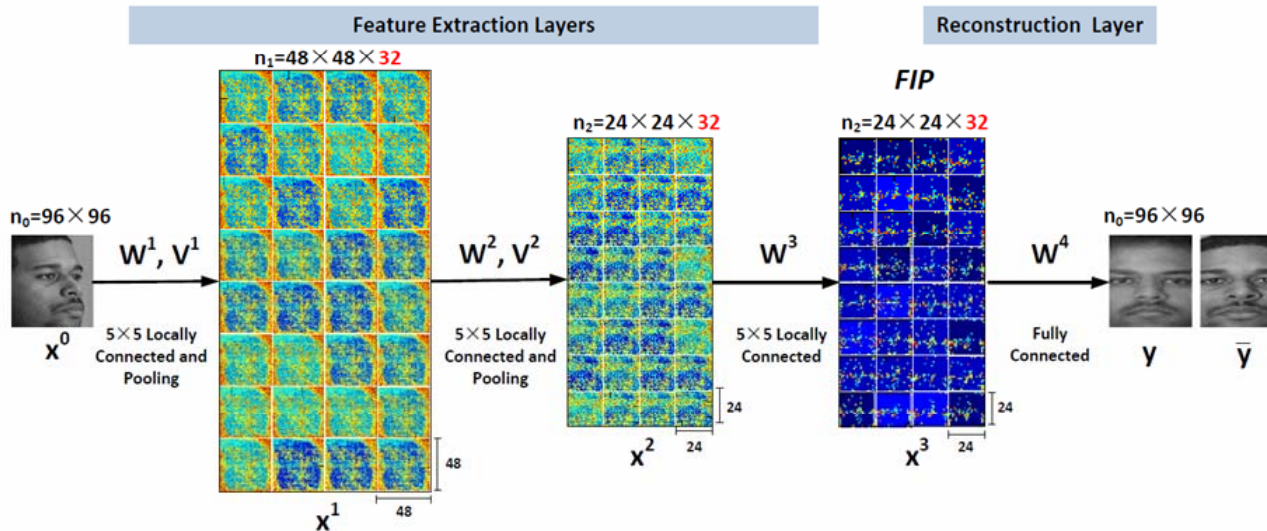
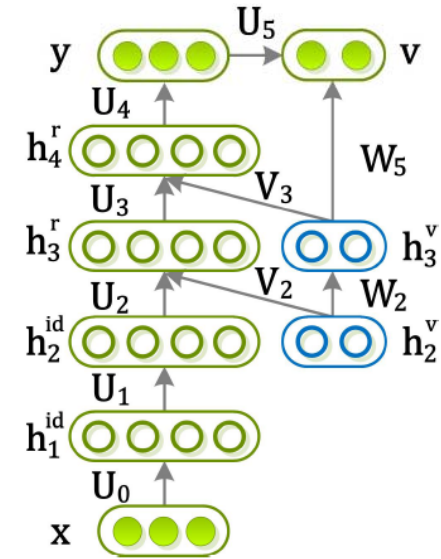
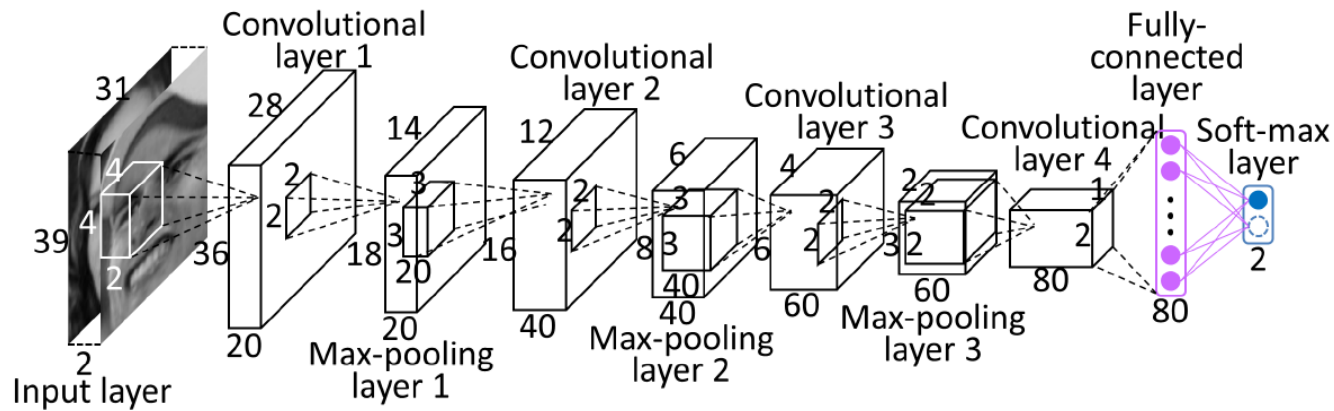


Ping Luo



Xiaoou Tang

Thank you!



<http://mmlab.ie.cuhk.edu.hk/>

<http://www.ee.cuhk.edu.hk/~xgwang/>