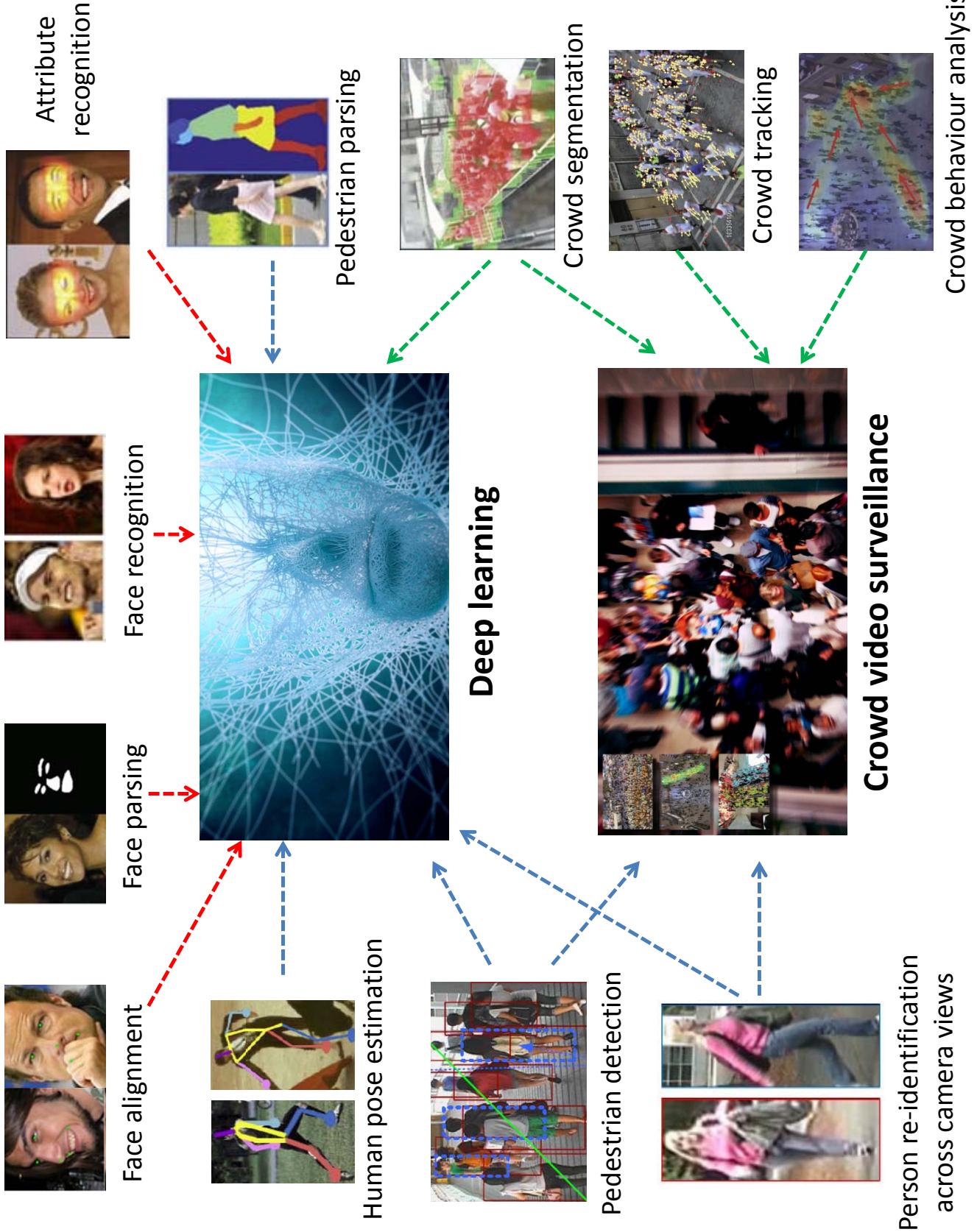


Deep Learning in Object Detection, Segmentation, and Recognition

Xiaogang Wang

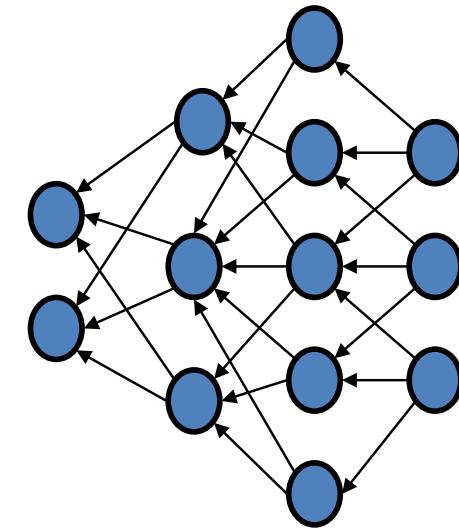
Department of Electronic Engineering,
The Chinese University of Hong Kong



Neural network
Back propagation

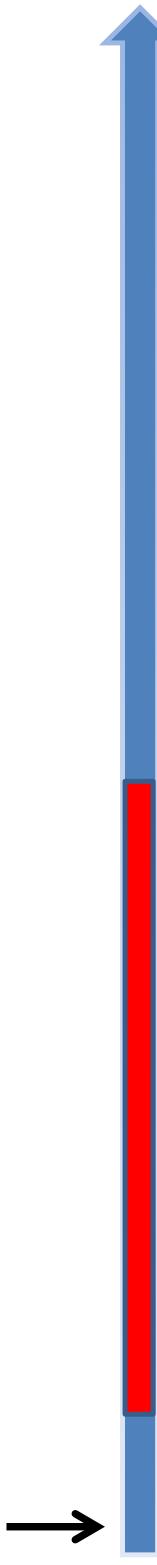


1986

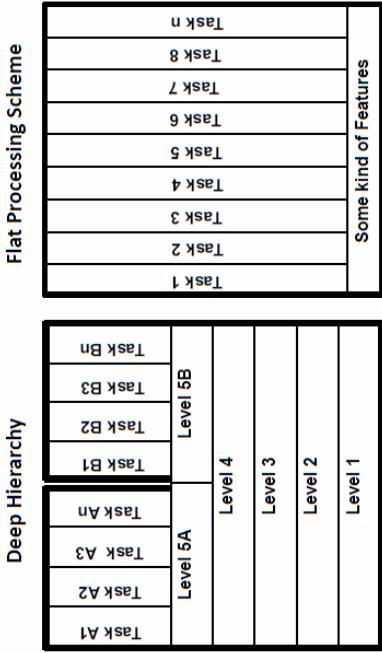


- Solve general learning problems
 - Tied with biological system
- But it is given up...
- Hard to train
 - Insufficient computational resources
 - Small training sets
 - Does not work well

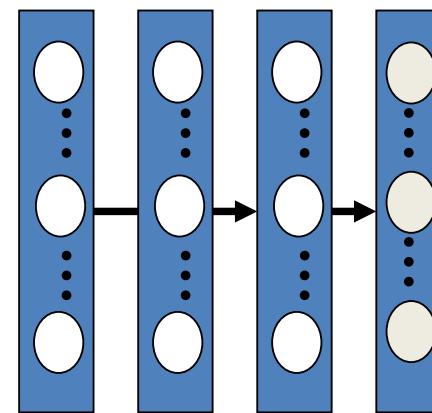
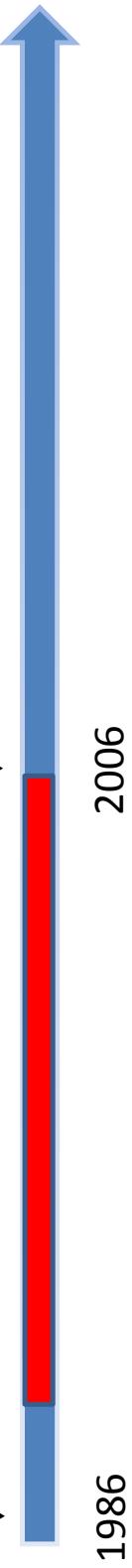
Neural network
Back propagation



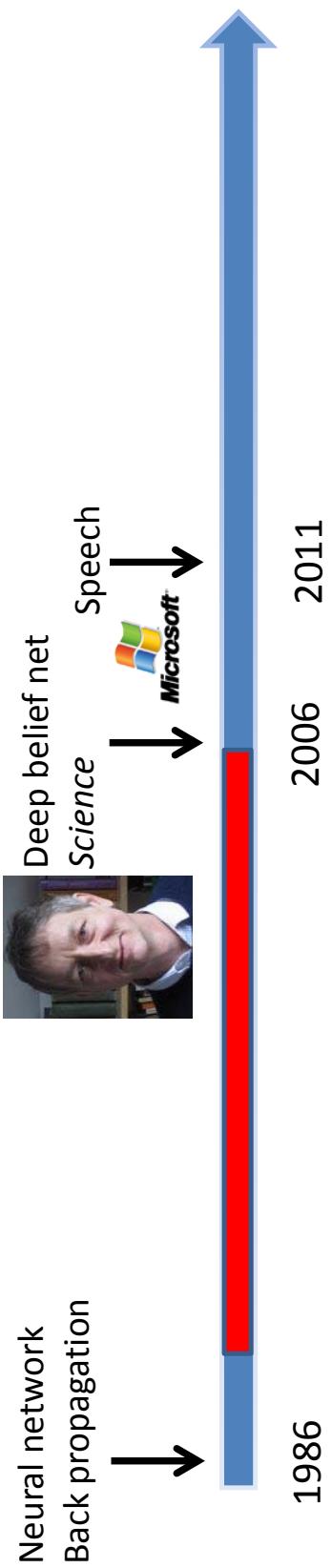
- SVM
- Boosting
- Decision tree
- KNN
- ...
- Loose tie with biological systems
- Flat structures
- Specific methods for specific tasks
 - Hand crafted features (GMM-HMM, SIFT, LBP, HOG)



Neural network
Back propagation



- Unsupervised & Layer-wised pre-training
- Better designs for modeling and training (normalization, nonlinearity, dropout)
- Feature learning
- New development of computer architectures
 - GPU
 - Multi-core computer systems
- Large scale databases



deep learning results

task	hours of training data	DNN-HMM	GMM-HMM with same data
Switchboard (test set 1)	309	18.5	27.4
Switchboard (test set 2)	309	16.1	23.6
English Broadcast News	50	17.5	18.8
Bing Voice Search (Sentence error rates)	24	30.4	36.2
Google Voice Input	5,870	12.3	
Youtube	1,400	47.6	52.3



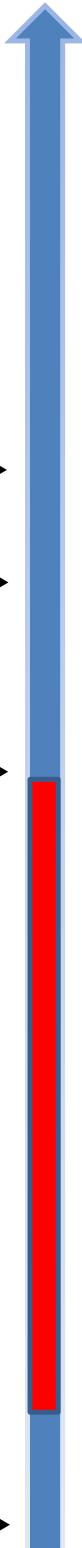
Deep Networks Advance State of Art in Speech

Deep Learning leads to breakthrough in speech recognition at MSR.

The New York Times



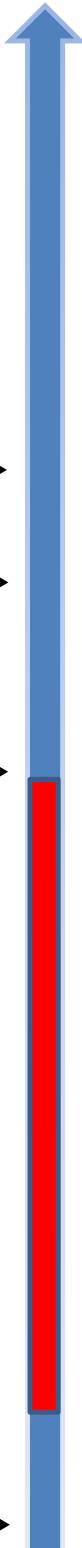
Neural network
Back propagation



1986

2006 2011 2012

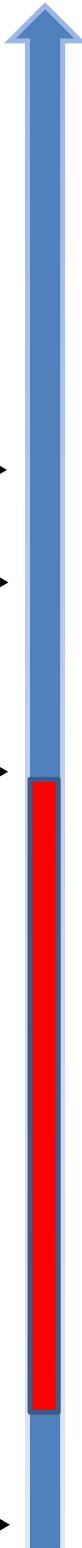
Deep belief net
Science



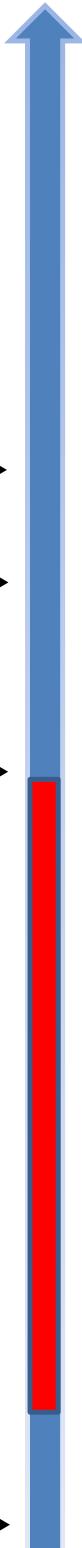
1986



Hong Kong



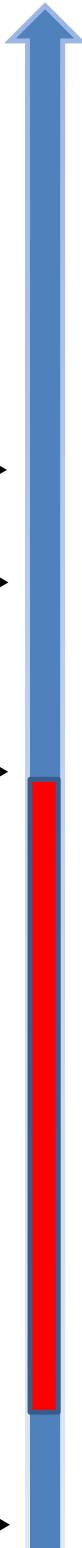
1986



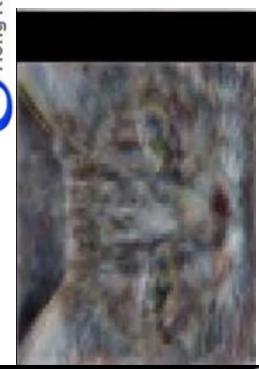
1986



Hong Kong



1986

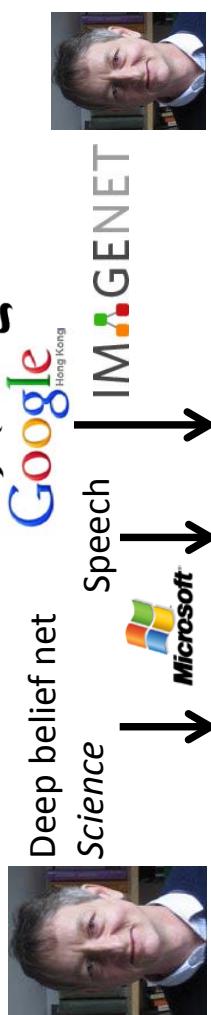


Trained on:

10,000,000 YouTube videos
1 frame from each (200x200)

How Many Computers to Identify a Cat? 16000 CPU cores

The New York Times

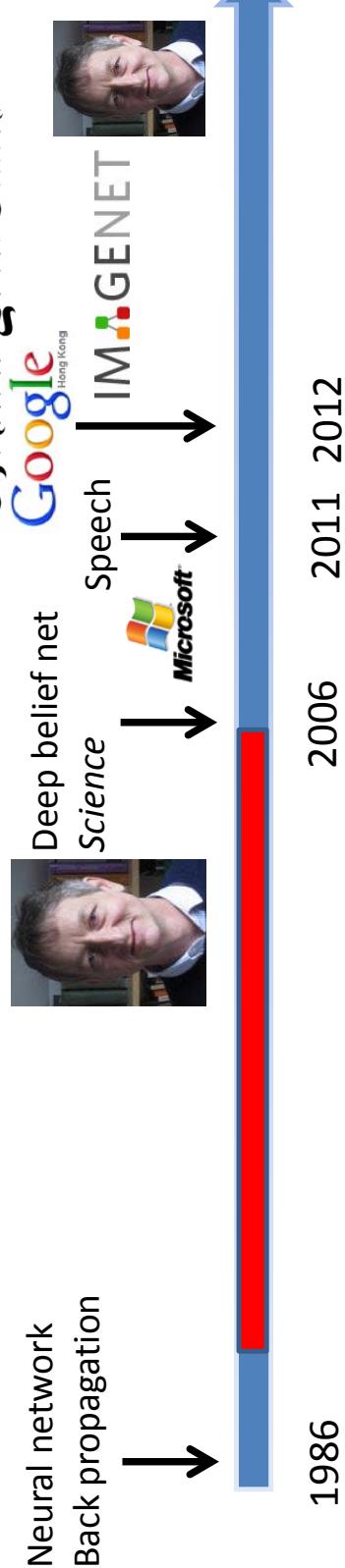


1986 2006 2011 2012

Rank	Name	Error rate	Description
1	U. Toronto	0.15315	Deep learning
2	U. Tokyo	0.26172	Hand-crafted
3	U. Oxford	0.26979	features and learning models.
4	Xerox/INRIA	0.27058	Bottleneck.

Object recognition over 1,000,000 images and 1,000 categories
(2 GPU)

The New York Times



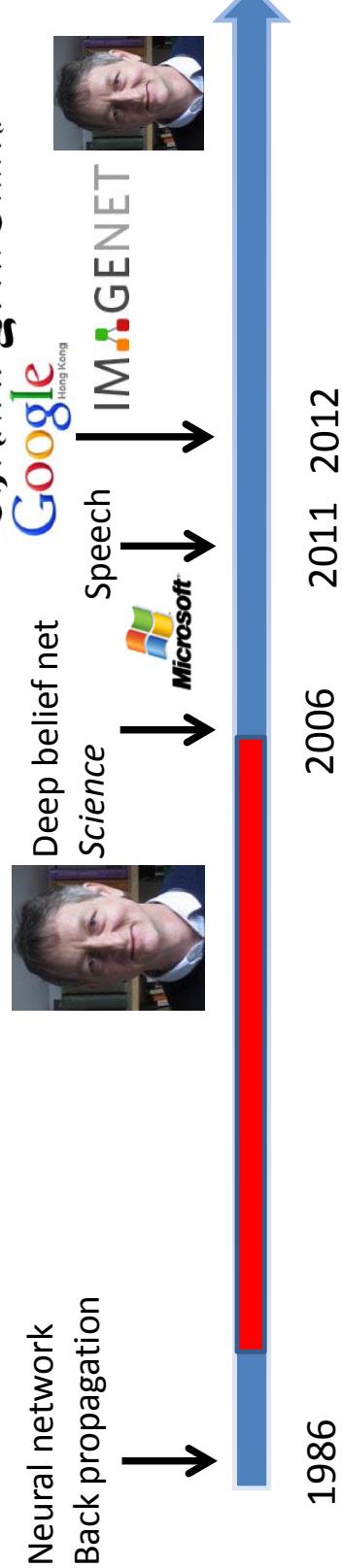
- ImageNet 2013

Rank	Name	Error rate	Description
1	NYU	0.11197	Deep learning
2	NUS	0.12535	Deep learning
3	Oxford	0.13555	Deep learning

MSRA, IBM, Adobe, NEC, Clarifai, Berkley, U. Tokyo, UCLA, UIUC, Toronto

Top 20 groups all used deep learning

The New York Times



- Google and Baidu announced their deep learning based visual search engines (2013)
 - Google
 - Baidu

Works Done by Us

Detection

- Pedestrian detection
- Facial keypoint detection

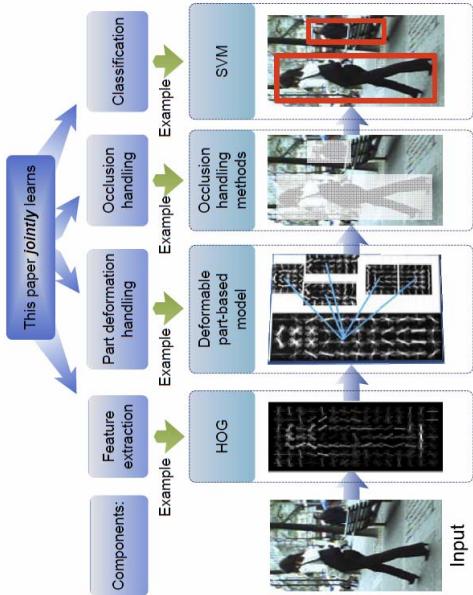
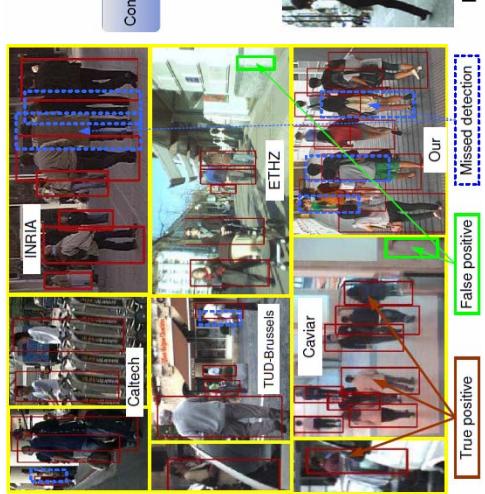
Segmentation

- Face parsing
- Pedestrian parsing

Recognition

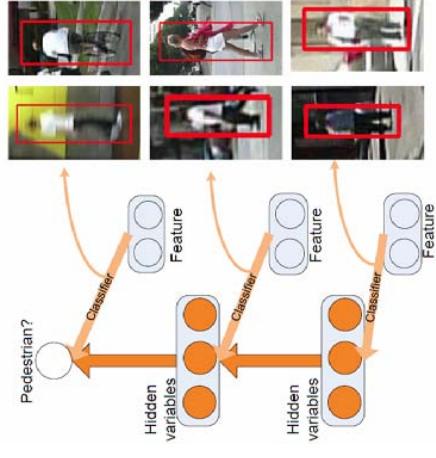
- Face verification
- Face attribute recognition

Pedestrian Detection

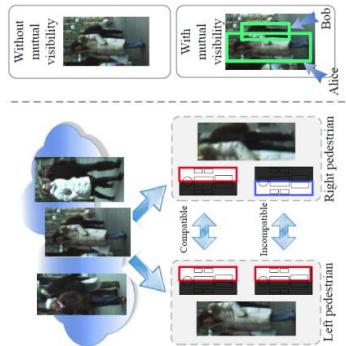


ICCV'13

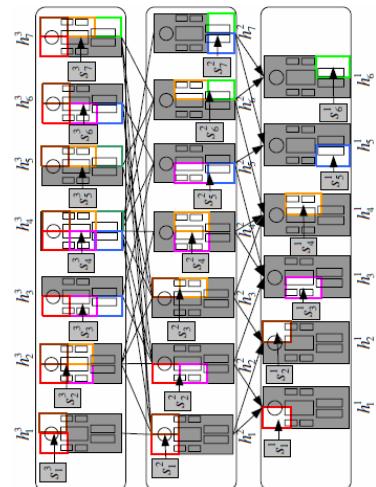
Improve state-of-the-art average miss detection rate on the largest Caltech dataset from 63% to 39%



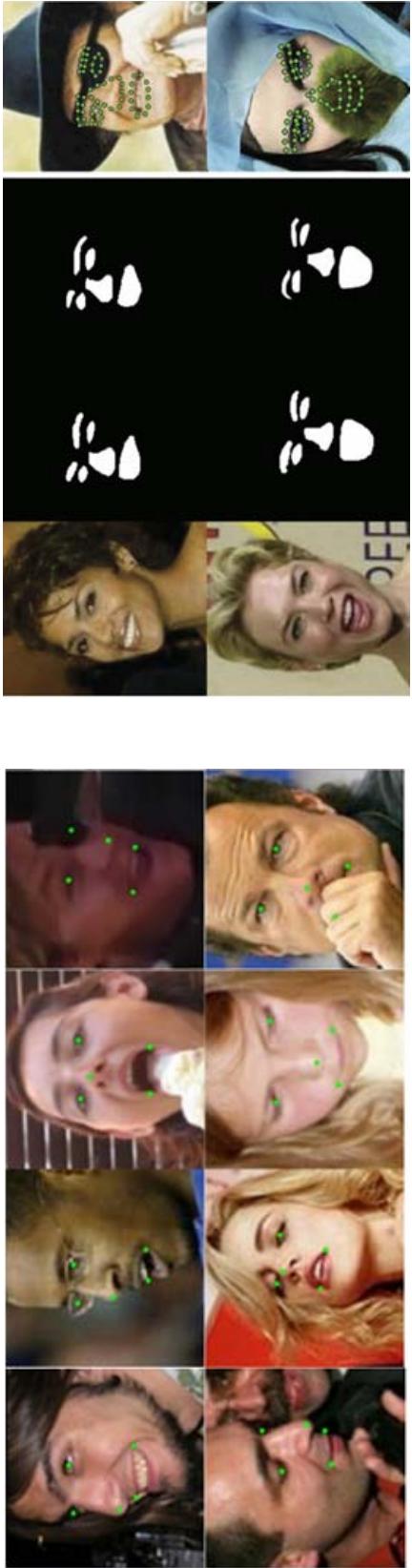
ICCV'13



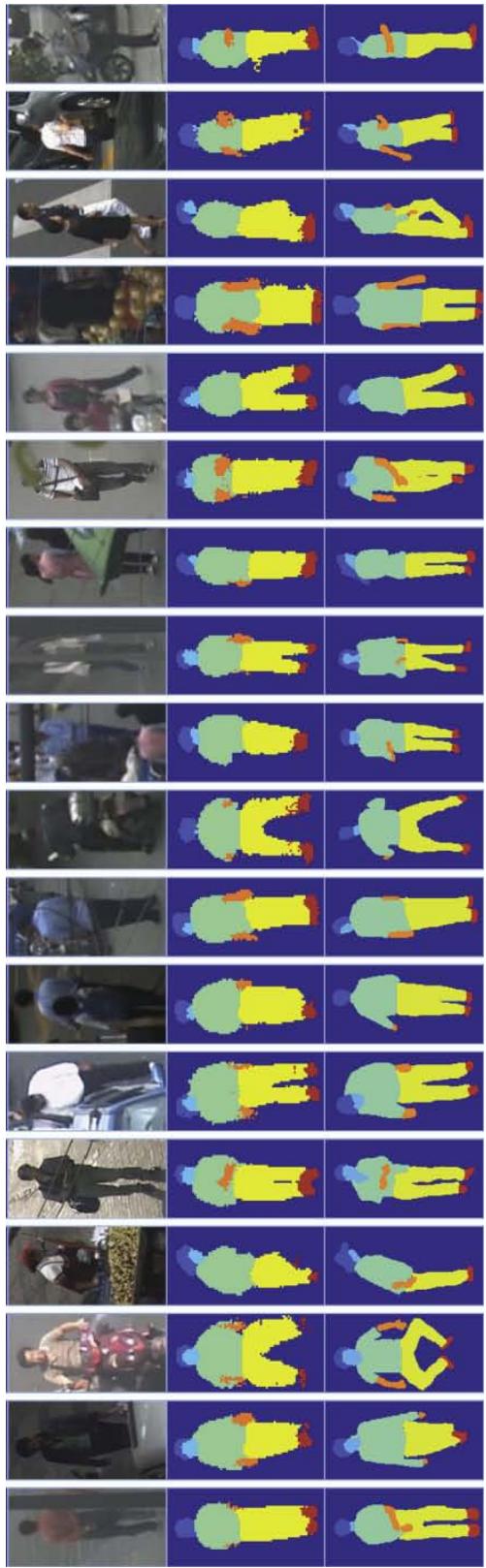
CVPR'13



CVPR'12



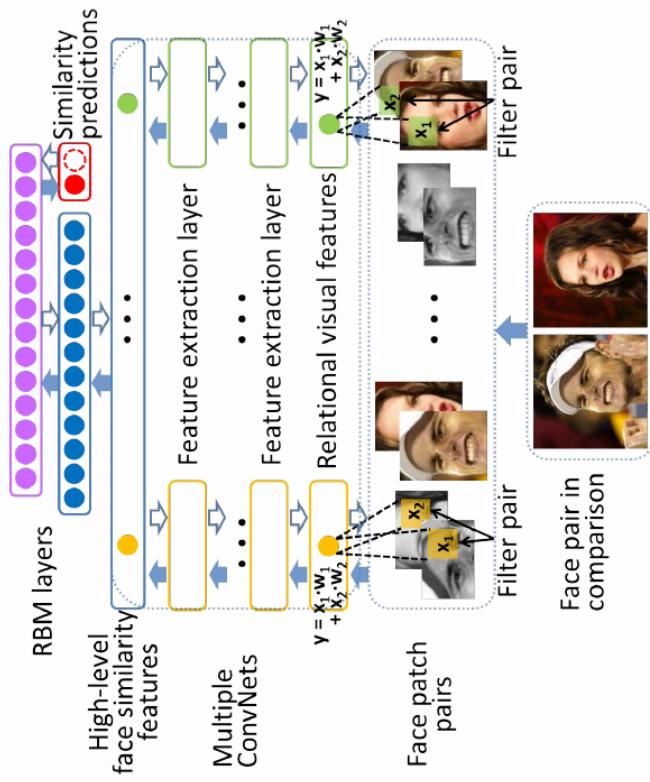
Facial keypoint detection, CVPR'13
(2% average error on LFWW)



Pedestrian parsing, CVPR'12

Face Recognition and Face Attribute Recognition

(LFW: 96.45%)



Face verification, ICCV'13

Recovering Canonical-View Face Images, ICCV'13



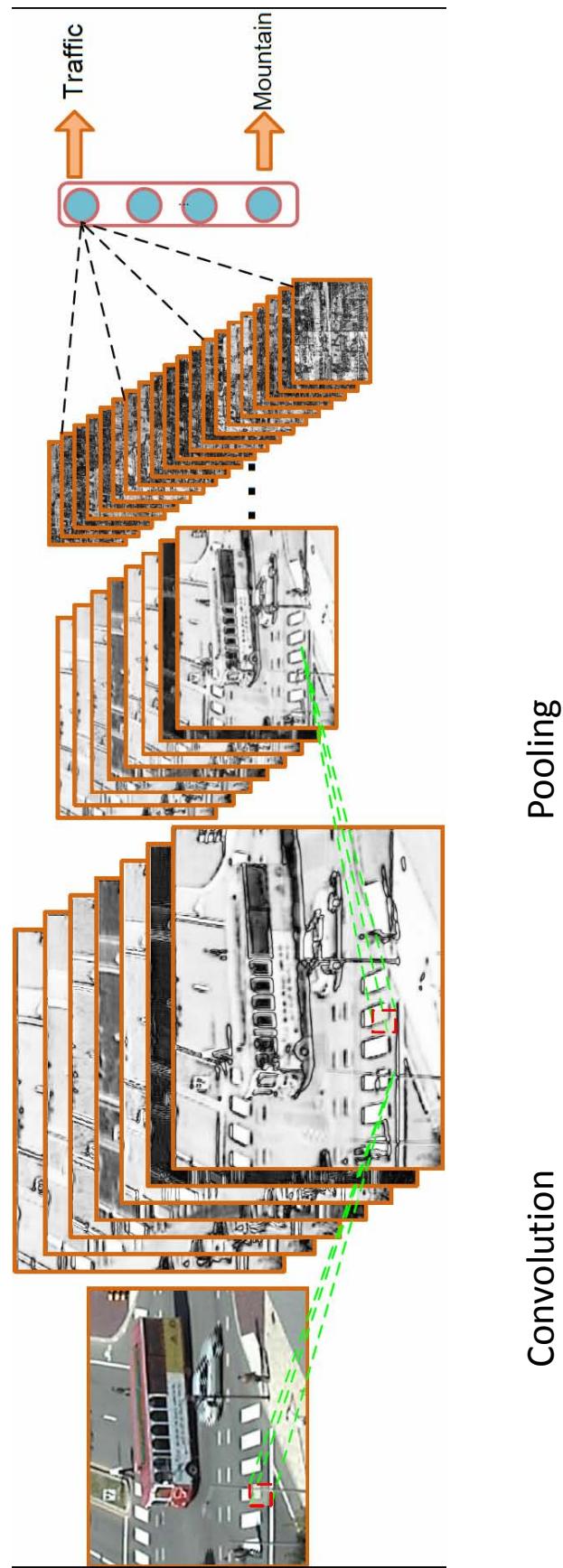
Face attribute recognition, ICCV'13

Introduction on Classical Deep Models

- Convolutional Neural Networks (CNN)
- Deep Belief Net (DBN)
- Auto-encoder

Classical Deep Models

- Convolutional Neural Networks (CNN)
 - LeCun'95



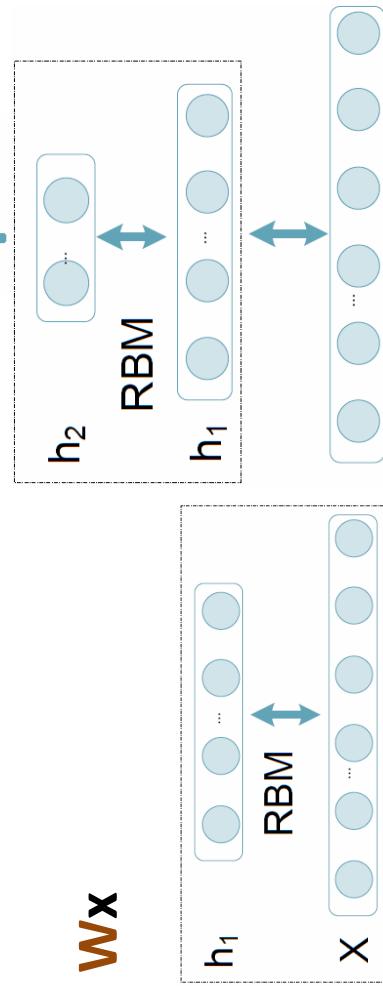
Classical Deep Models

- Deep belief net
– Hinton'06

$$P(\mathbf{x}, \mathbf{h}_1, \mathbf{h}_2) = p(\mathbf{x} | \mathbf{h}_1) p(\mathbf{h}_1, \mathbf{h}_2)$$

$$P(\mathbf{x}, \mathbf{h}_1) = \frac{e^{-E(\mathbf{x}, \mathbf{h}_1)}}{\sum_{\mathbf{x}, \mathbf{h}_1} e^{-E(\mathbf{x}, \mathbf{h}_1)}}$$

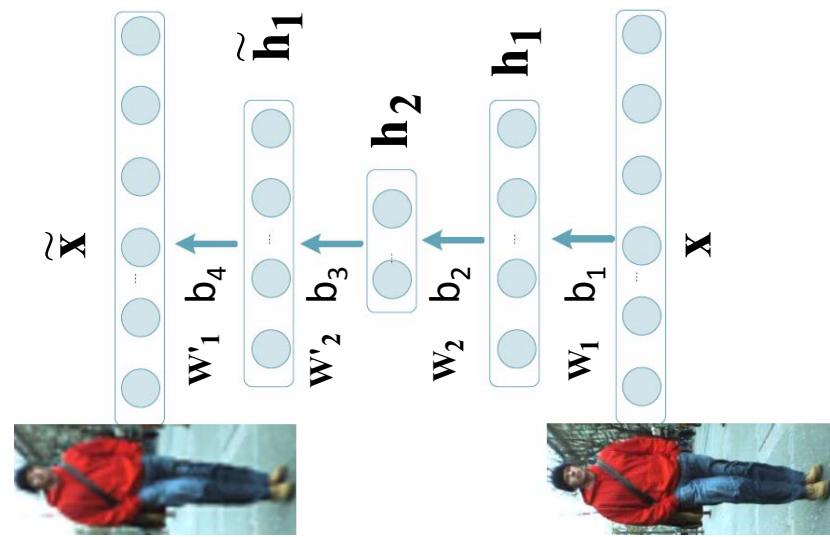
$$E(\mathbf{x}, \mathbf{h}_1) = \mathbf{b}' \mathbf{x} + \mathbf{c}' \mathbf{h}_1 + \mathbf{h}_1' \mathbf{W} \mathbf{x}$$



Classical Deep Models

- Auto-encoder

– Hinton'06

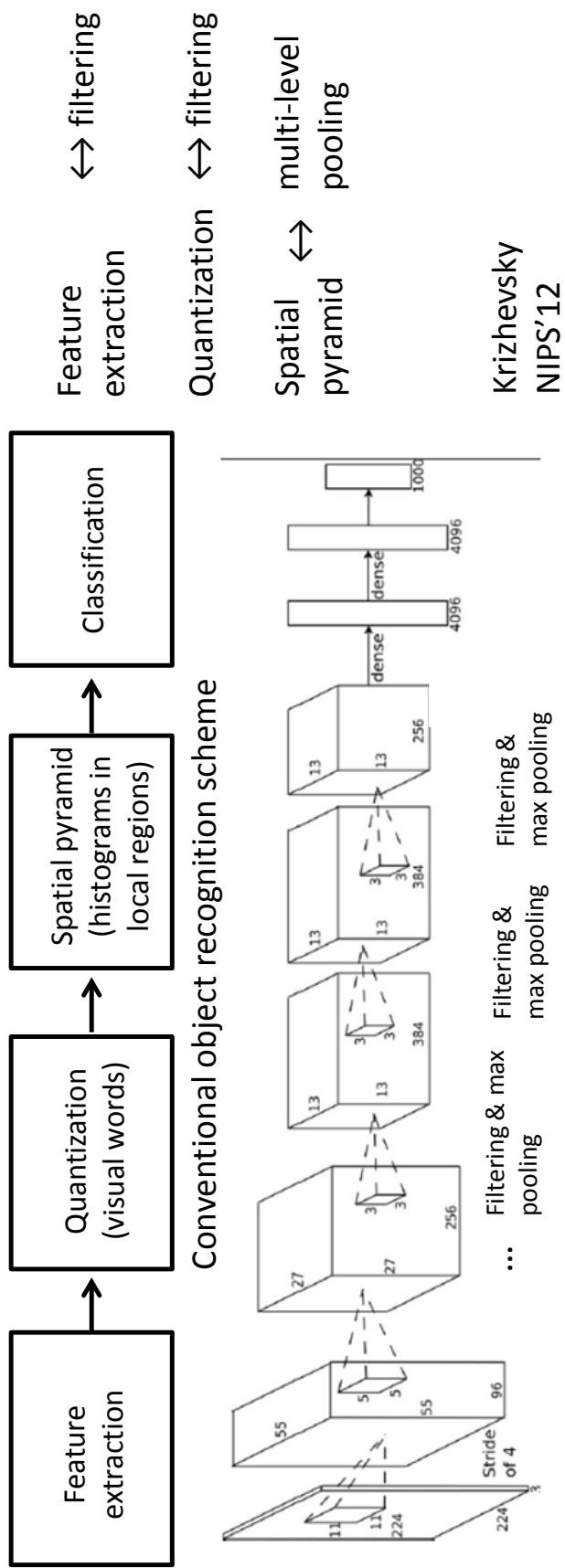


$$\begin{aligned}\text{Encoding: } \mathbf{h}_1 &= \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \\ \mathbf{h}_2 &= \sigma(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2)\end{aligned}$$

$$\begin{aligned}\text{Decoding: } \tilde{\mathbf{h}}_1 &= \sigma(\mathbf{W}'_2 \mathbf{h}_2 + \mathbf{b}_3) \\ \tilde{\mathbf{x}} &= \sigma(\mathbf{W}'_1 \mathbf{h}_1 + \mathbf{b}_4)\end{aligned}$$

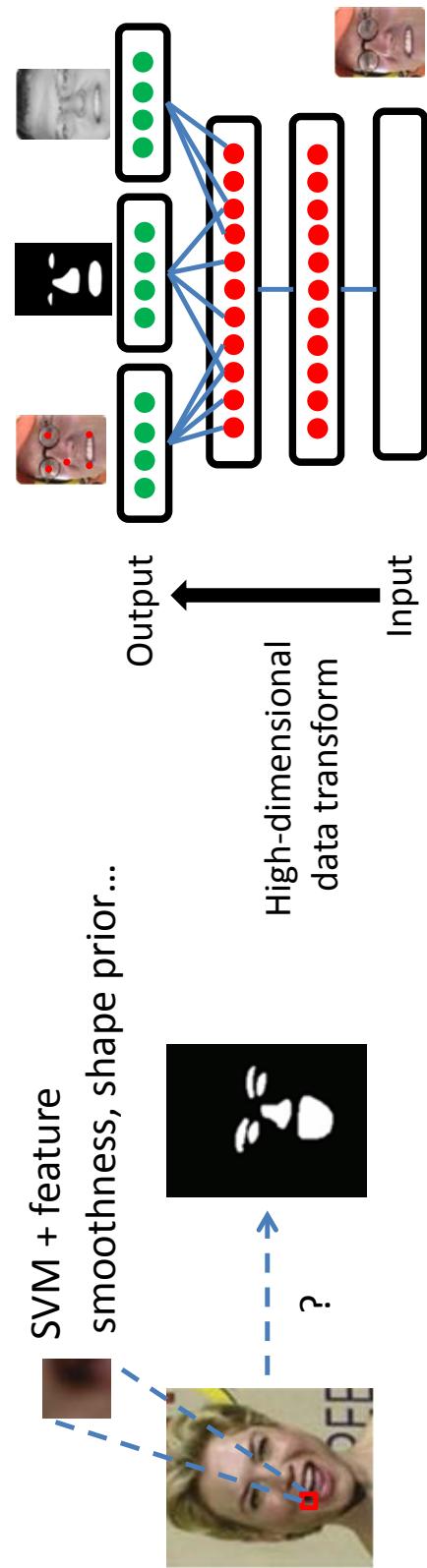
Opinion I

- How to formulate a vision problem with deep learning?
 - Make use of experience and insights obtained in CV research
 - Sequential design/learning vs **joint learning**
 - Effectively train a deep model (layerwise pre-training + fine tuning)



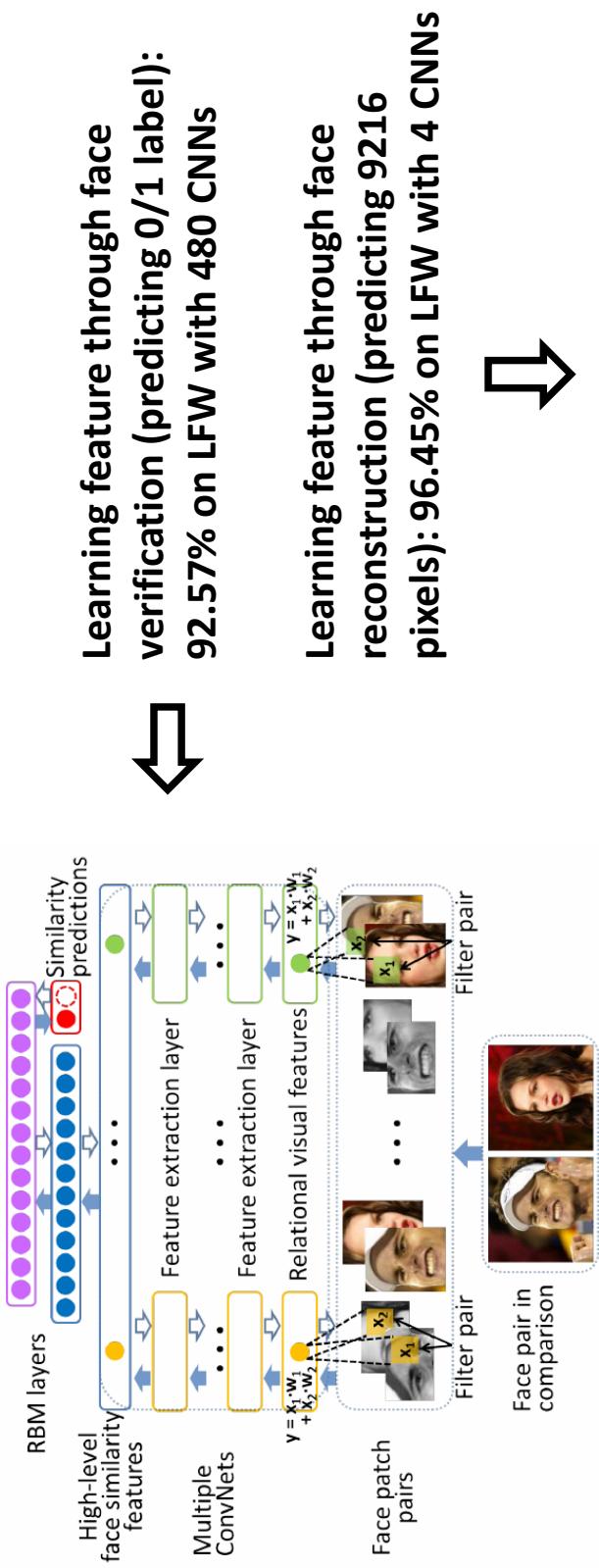
Opinion II

- How to make use of the large learning capacity of deep models?
 - **High dimensional data transform**
 - Hierarchical nonlinear representations

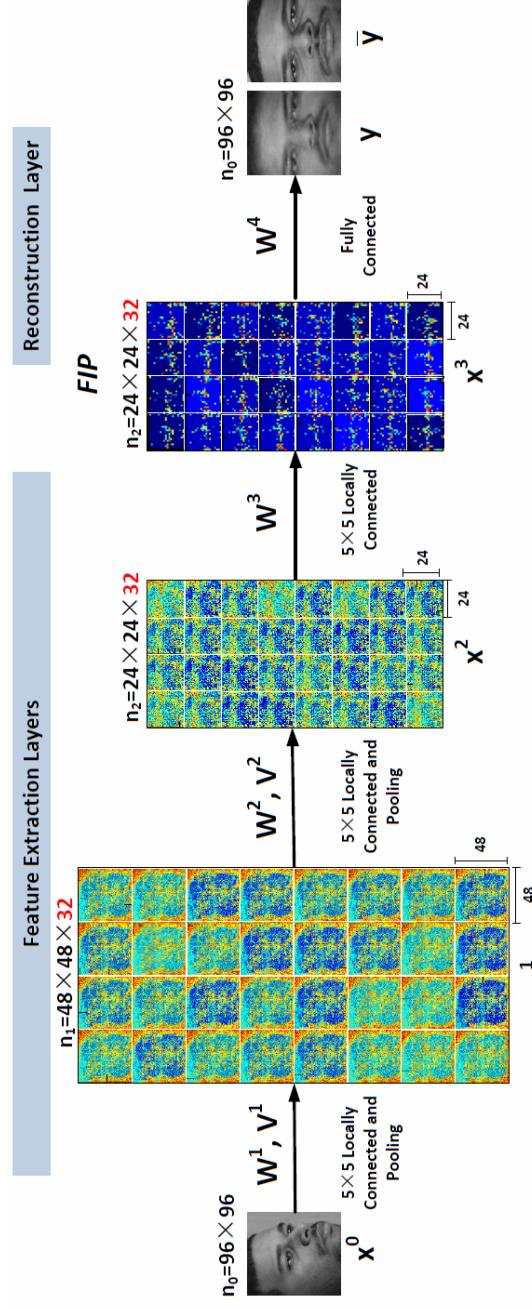


Opinion III

- Deep learning likes challenging tasks (for better generalization)
 - Make input data more challenging (augmenting data by translating, rotating, and scaling)
 - Make training process more challenging (dropout: randomly setting some responses to zero; dropconnect: randomly setting some weights to zero)
 - **Make prediction more challenging**



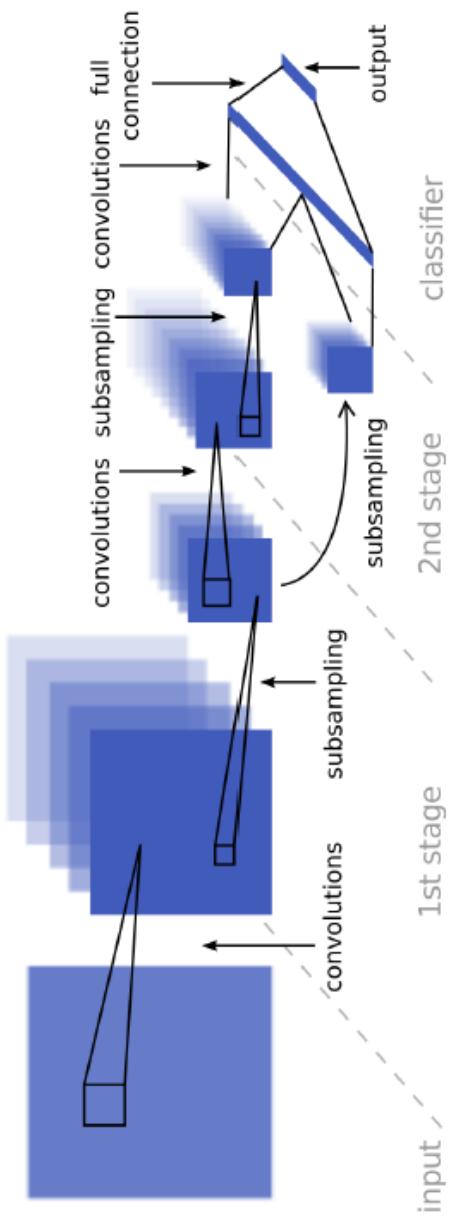
Y. Sun, X. Wang, and X. Tang, "Hybrid Deep Learning for Computing Face Similarities," ICCV'13



Z. Zhu, P. Luo, X. Wang, and X. Tang, "Deep Learning Identify-Preserving Face Space," ICCV 2013.

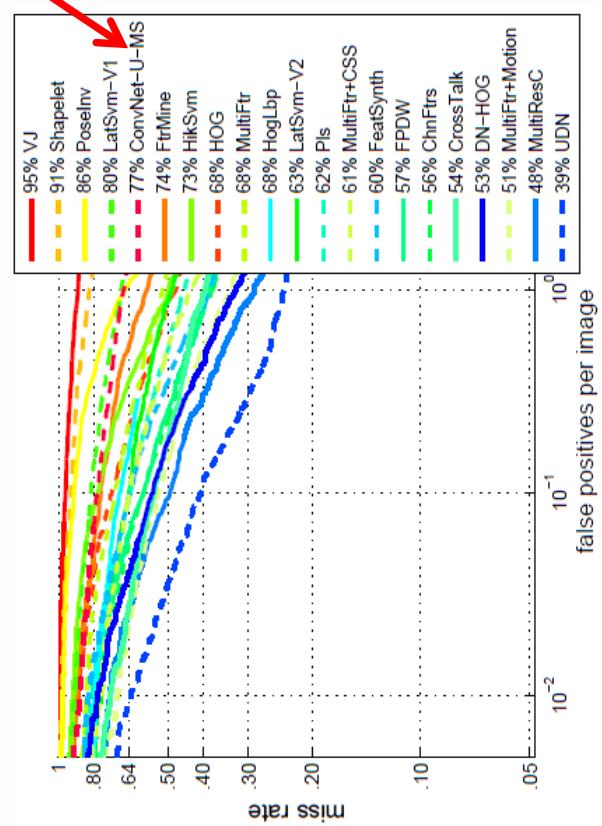
Joint Deep Learning

What if we treat an existing deep model as a black box in pedestrian detection?

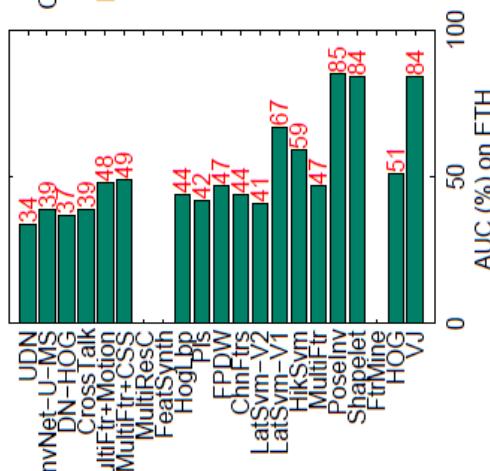


ConvNet-U-MS

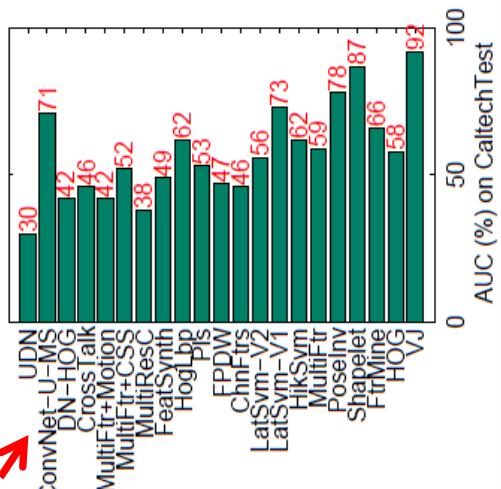
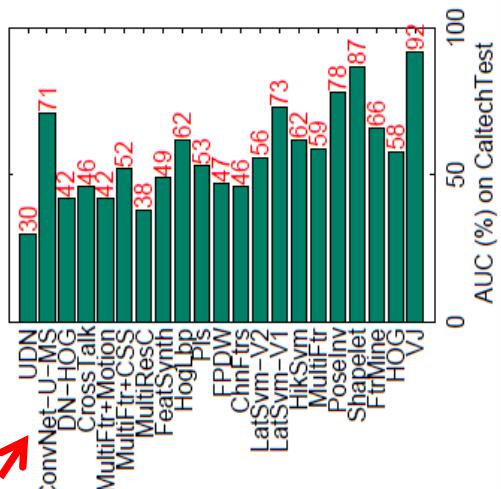
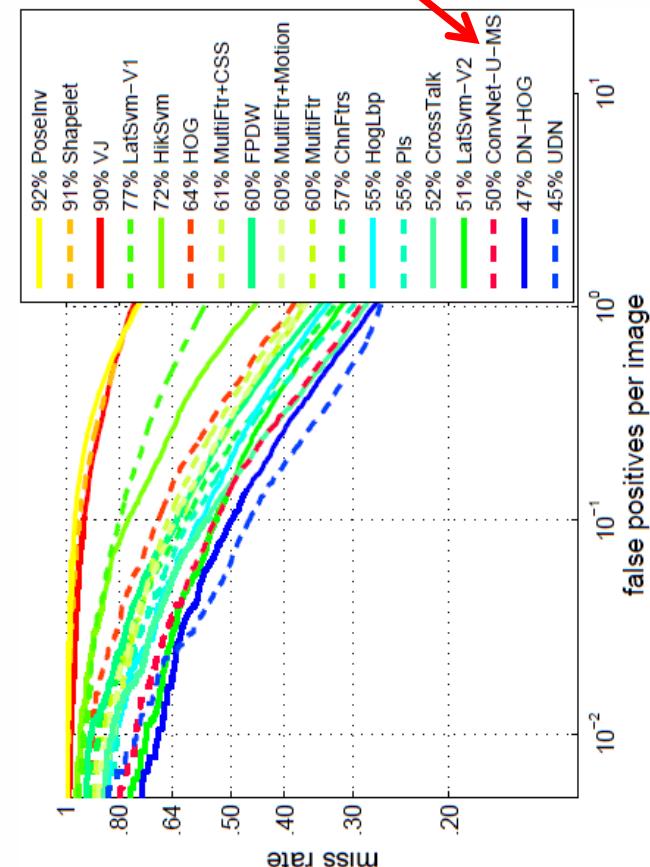
- Sermnet, K. Kavukcuoglu, S. Chintala, and LeCun, “Pedestrian Detection with Unsupervised Multi-Stage Feature Learning,” CVPR 2013.

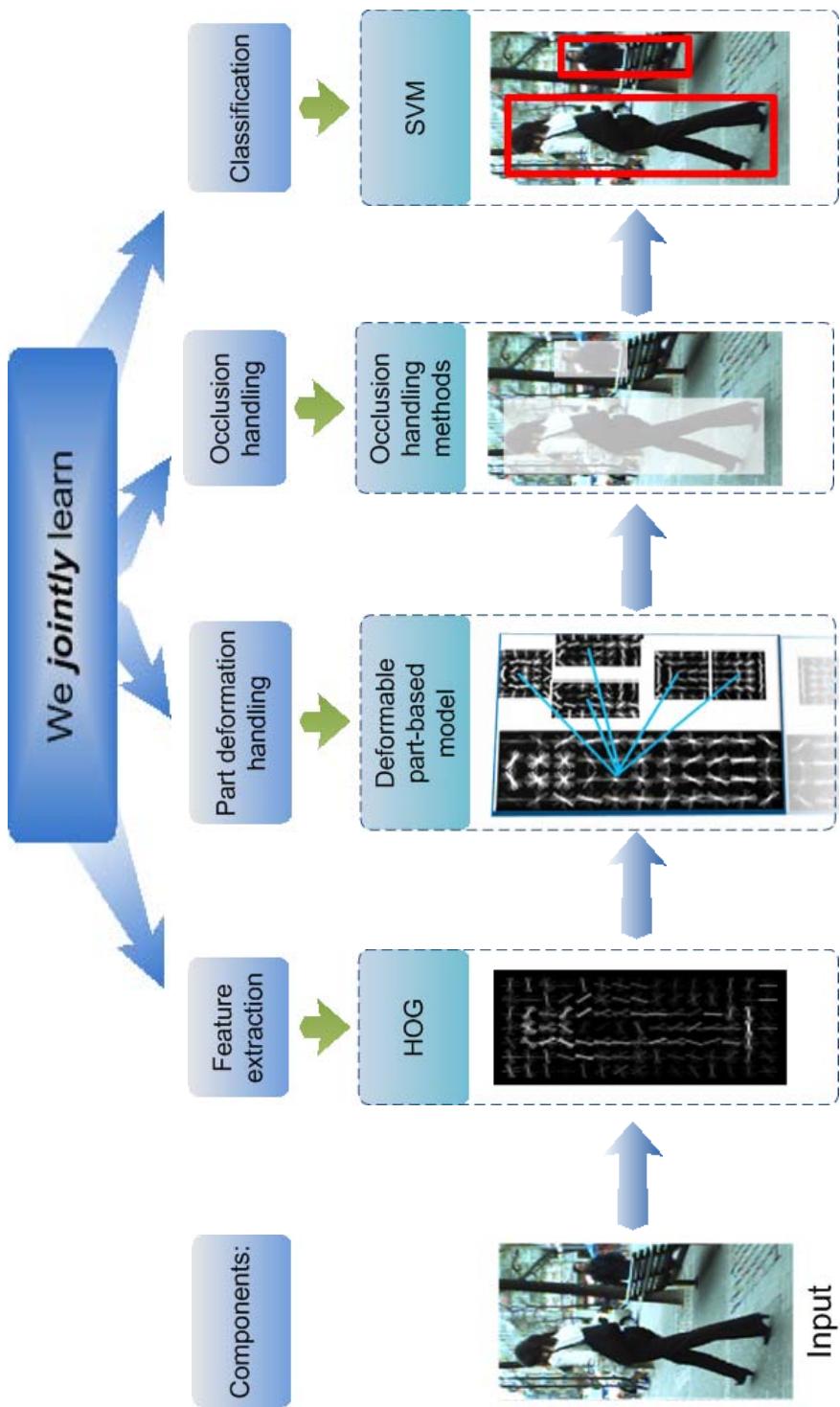


Results on Caltech Test



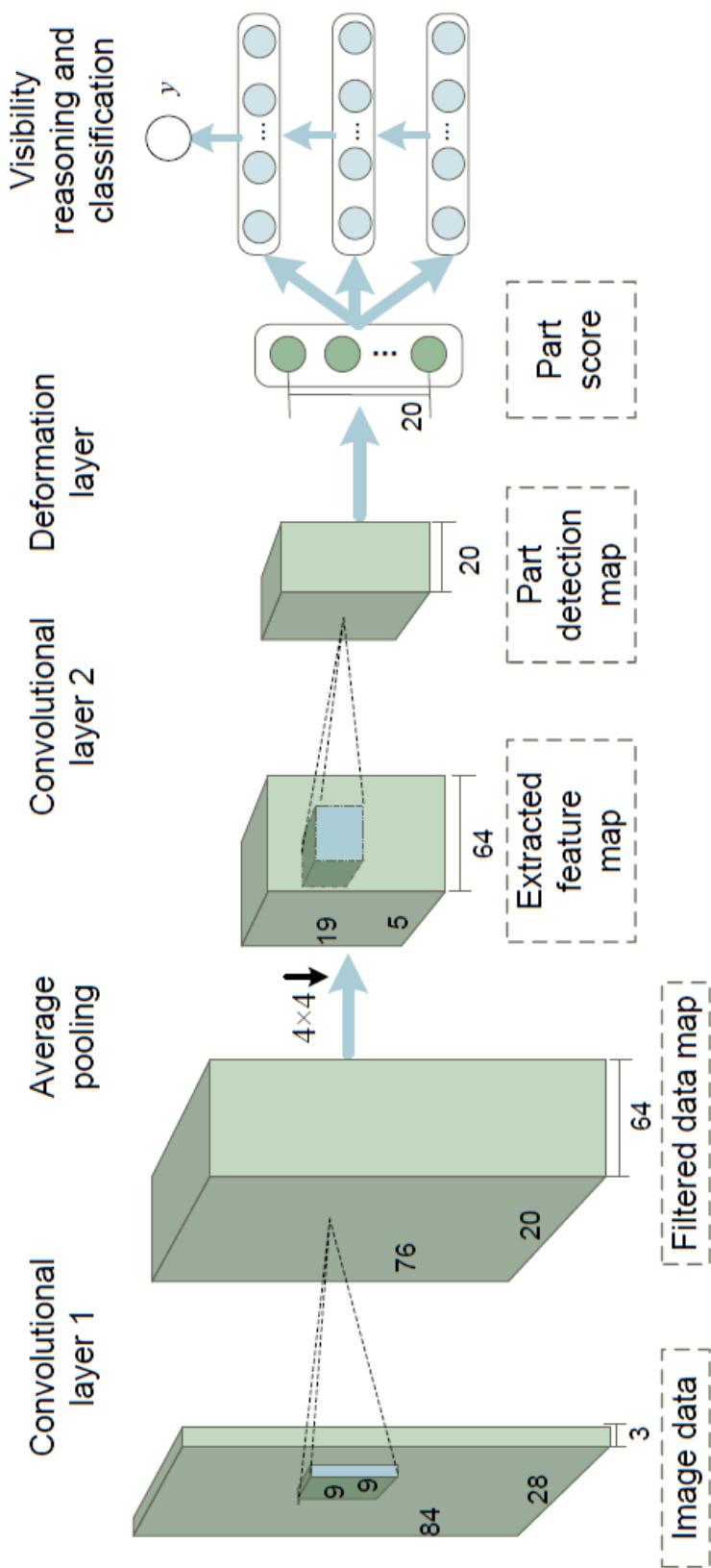
Results on ETHZ





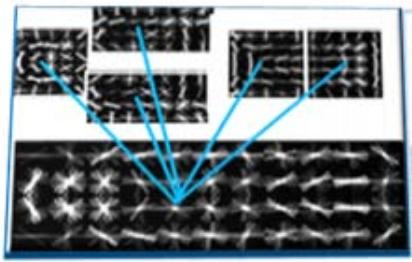
- N. Dalal and B. Triggs. Histograms of oriented gradients for human detection.
- P. Felzenszwalb, D. McAllester, and D. Ramanan. A Discriminatively Trained, Multiscale, Deformable Part Model. CVPR, 2008. (2000 citations)
- W. Ouyang and X. Wang. A Discriminative Deep Model for Pedestrian Detection with Occlusion Handling. CVPR, 2012.

Our Joint Deep Learning Model

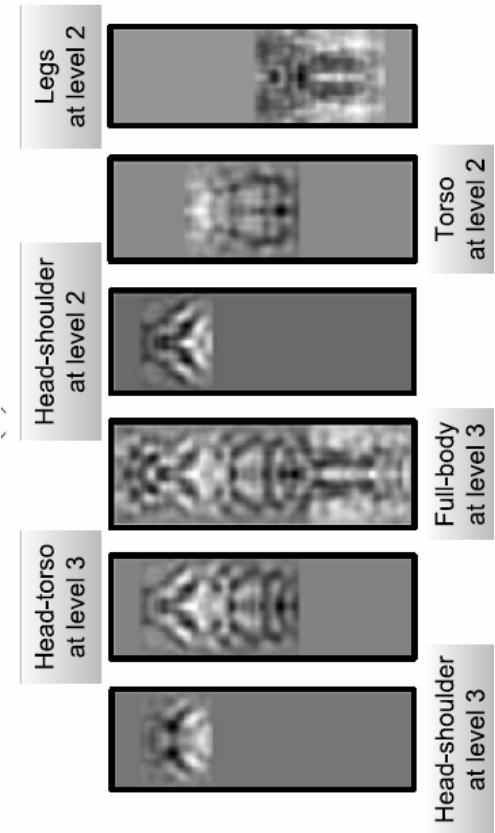


Modeling Part Detectors

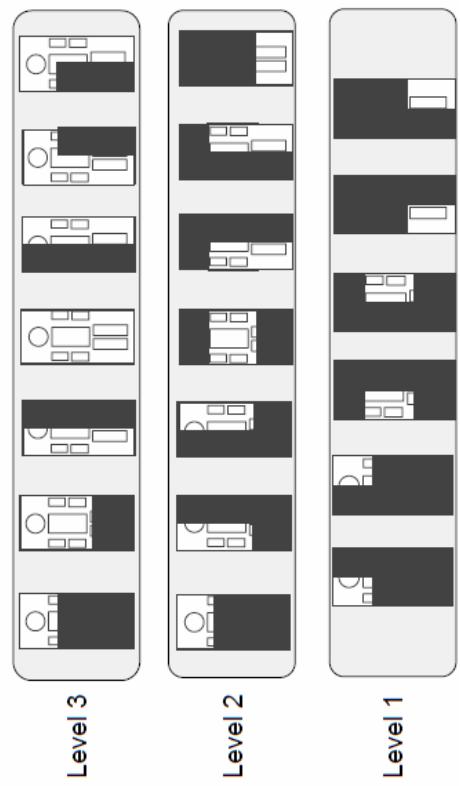
- Design the filters in the second convolutional layer with variable sizes



Part models learned
from HOG

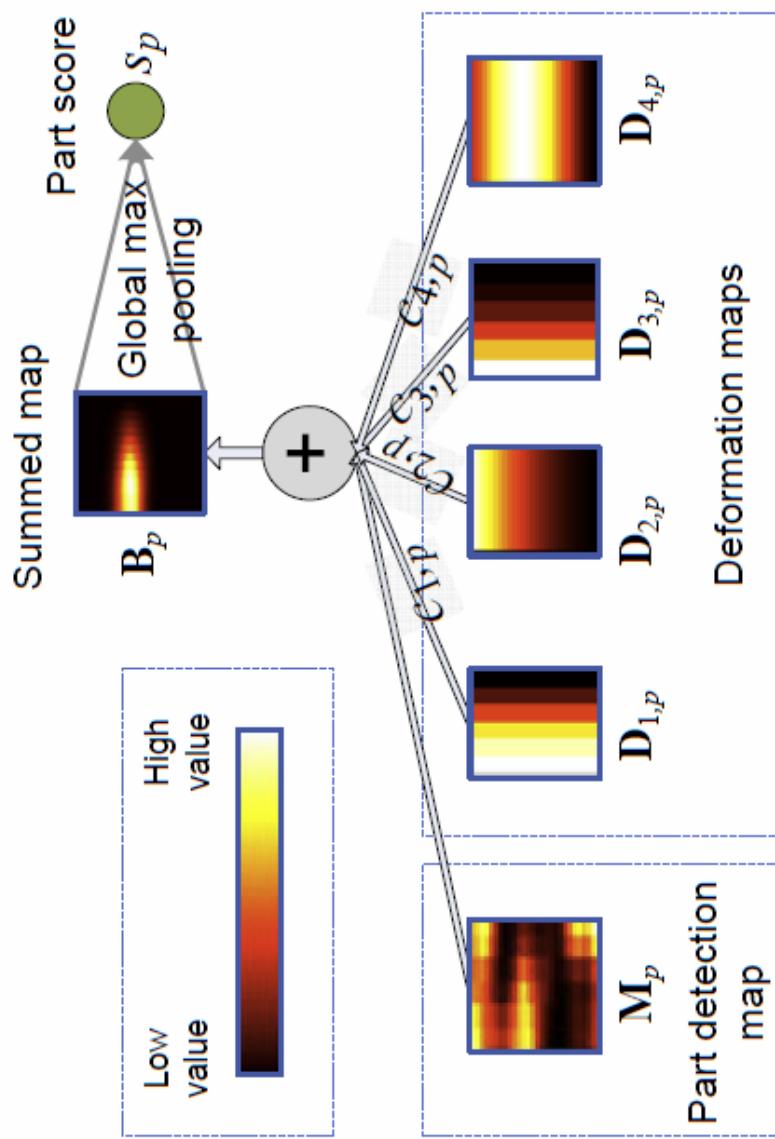


Learned filtered at the second
convolutional layer

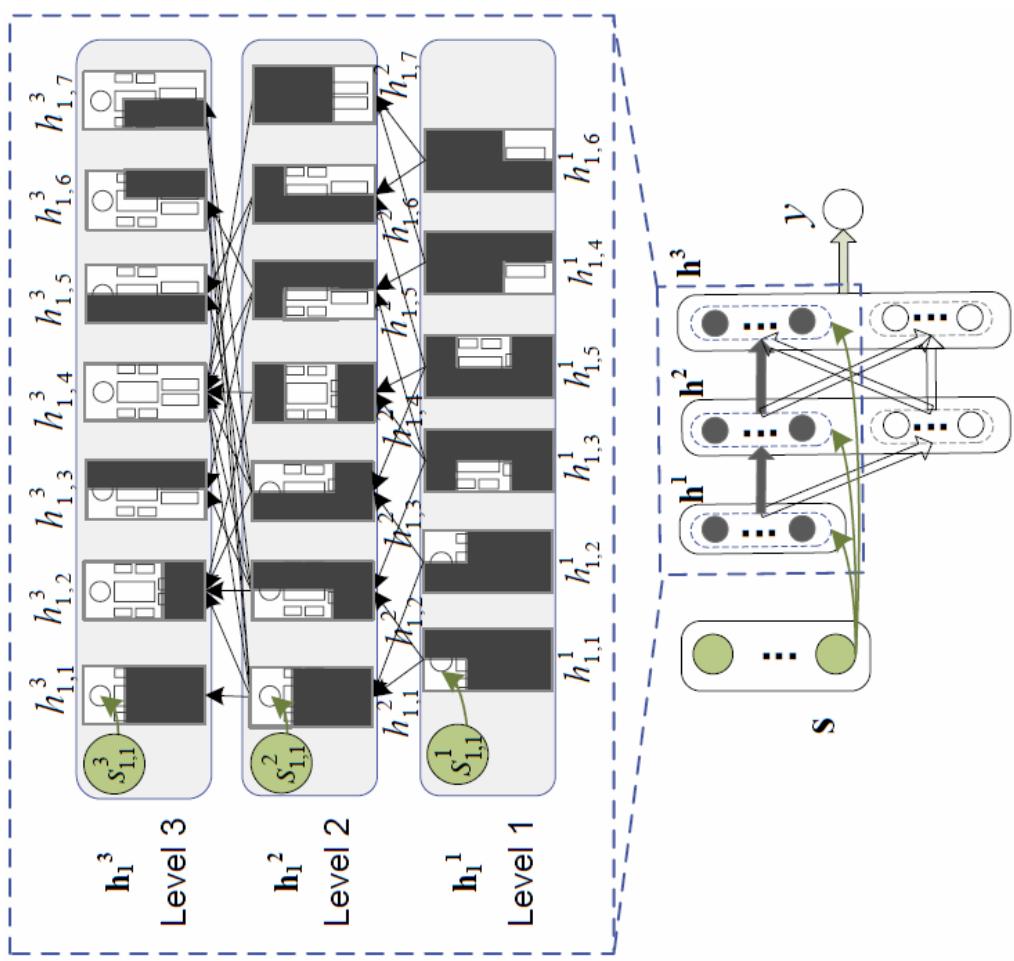


Part models

Deformation Layer

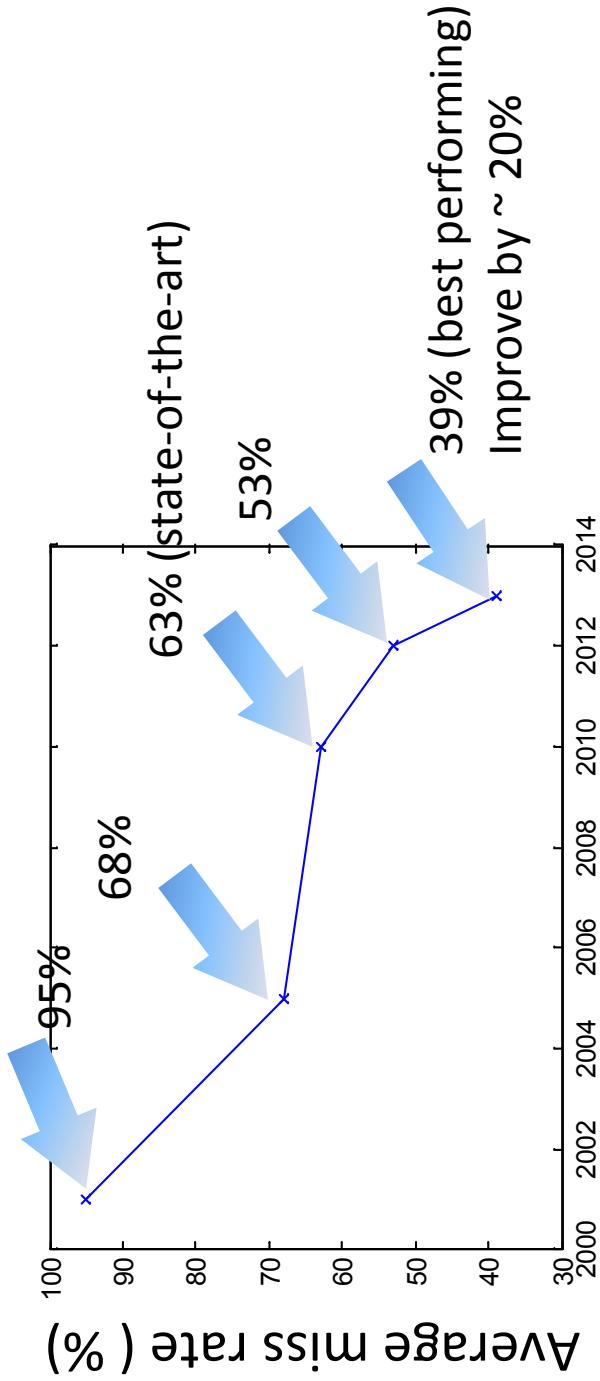


Visibility Reasoning with Deep Belief Net



Experimental Results

- Caltech – Test dataset (largest, most widely used)



W. Ouyang and X. Wang, "A Discriminative Deep Model for Pedestrian Detection with Occlusion Handling," CVPR 2012.

W. Ouyang, X. Zeng, and X. Wang, "Multi-Pedestrian Detection aided by Pedestrian Detection Model," CVPR 2013.

W. Ouyang, X. Zeng, and X. Wang, "Single Pedestrian Detection aided by Multi-Pedestrian Detection," CVPR 2013.

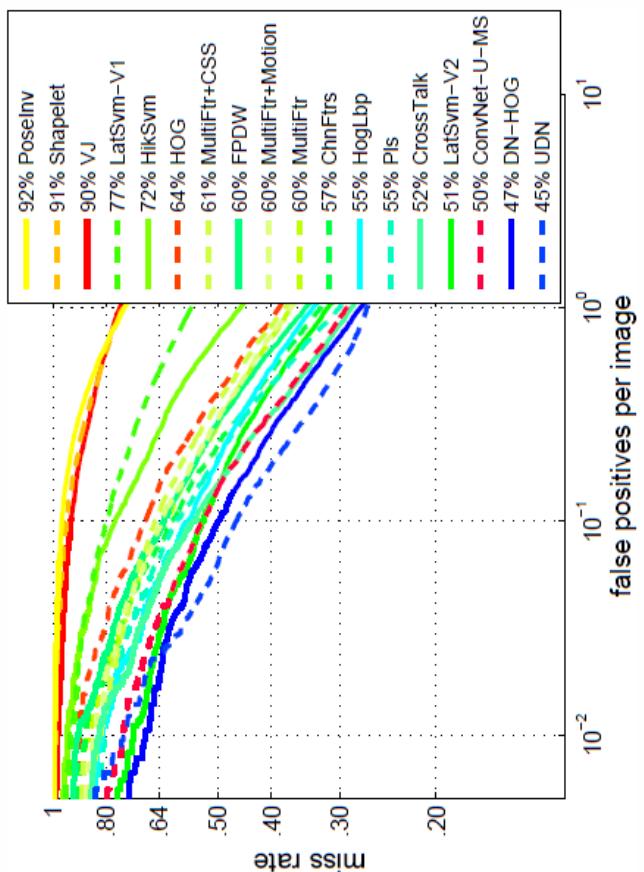
X. Zeng, W. Ouyang, and X. Wang, "A Cascaded Dense Learning Architecture for Pedestrian Detection," CVPR 2013.

W. Ouyang, X. Zeng, and X. Wang, "Joint Deep Learning for Pedestrian Detection," IEEE Conference on Computer Vision and Pattern Recognition (CVPR) 2013.

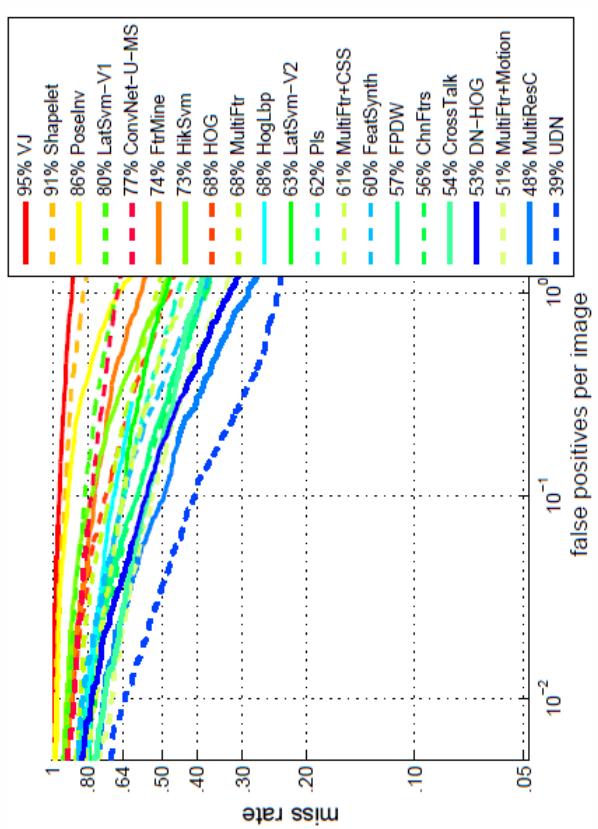
This work is distinguished by three key contributions. The first is the introduction of a new ...

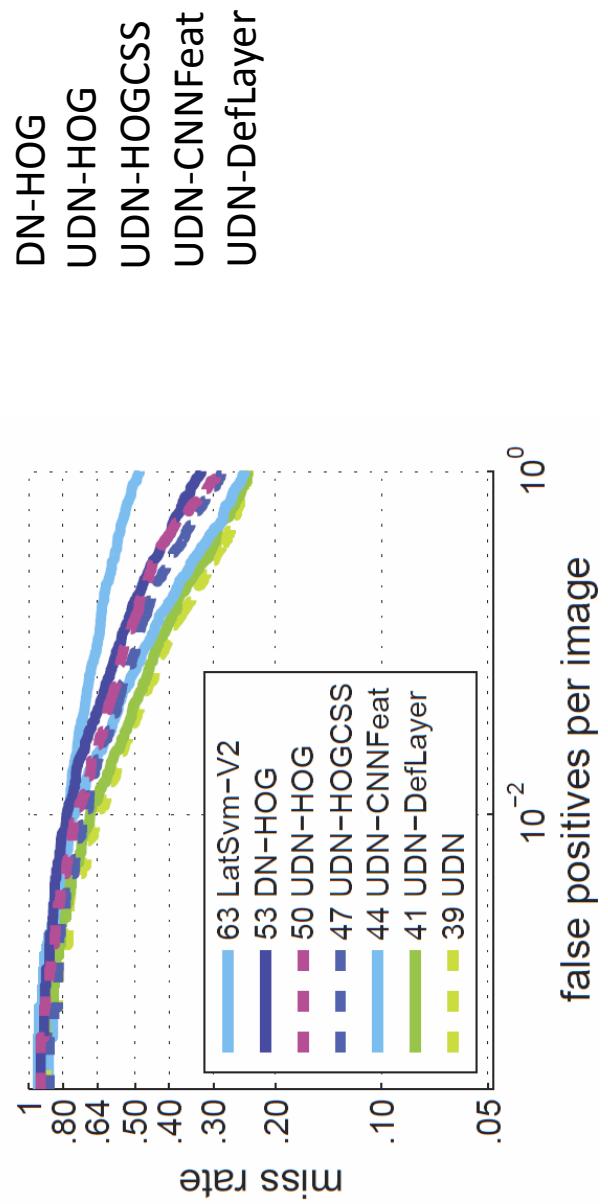
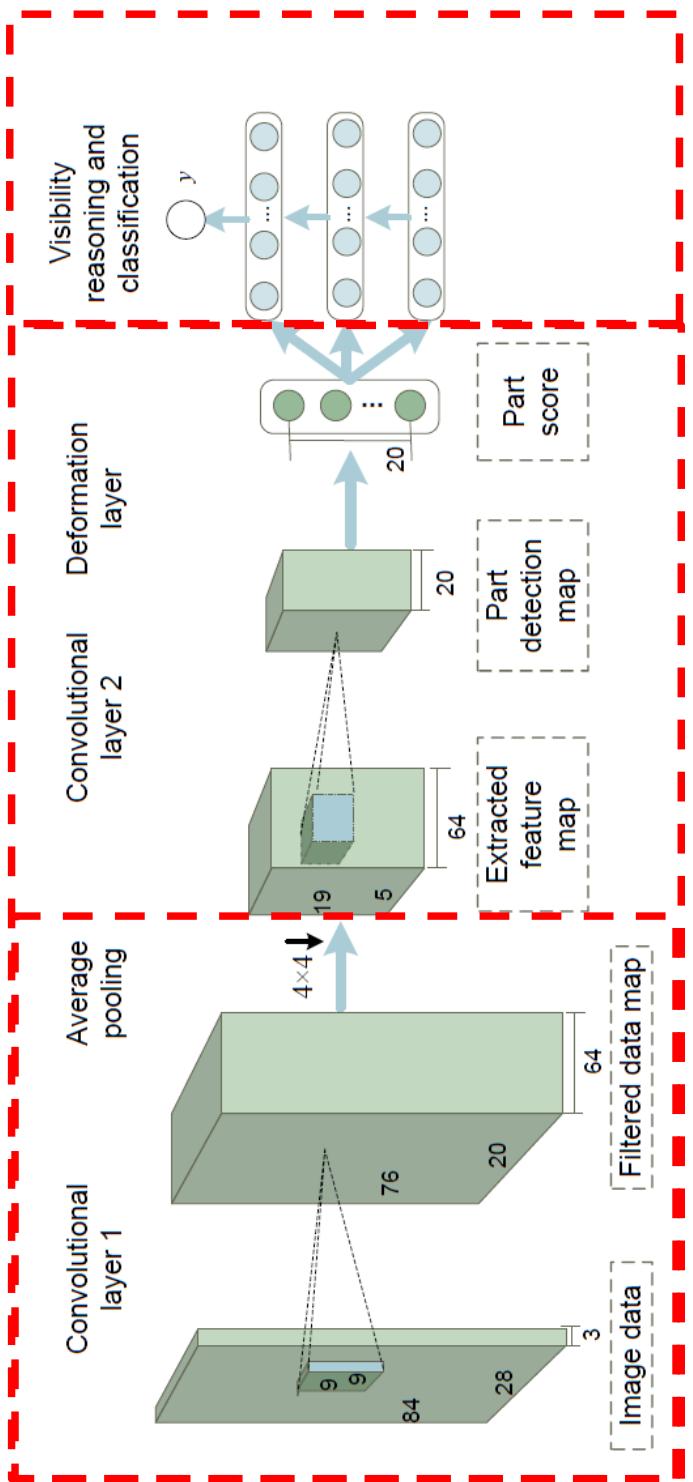
Cited by 7647 Related articles All 2011 versions Import into BibTeX More ▾

Results on ETHZ



Results on Caltech Test

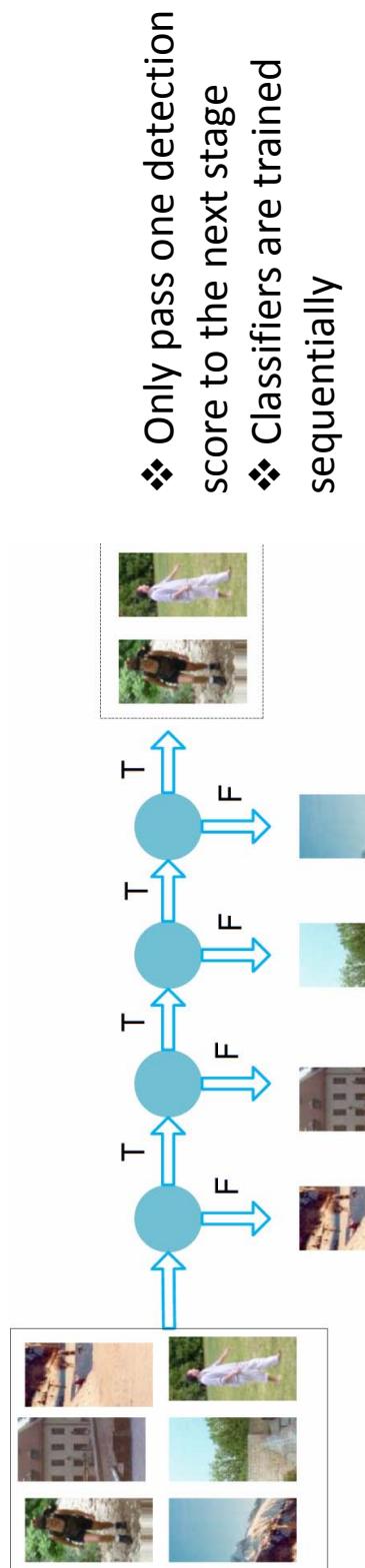




Multi-Stage Contextual Deep Learning

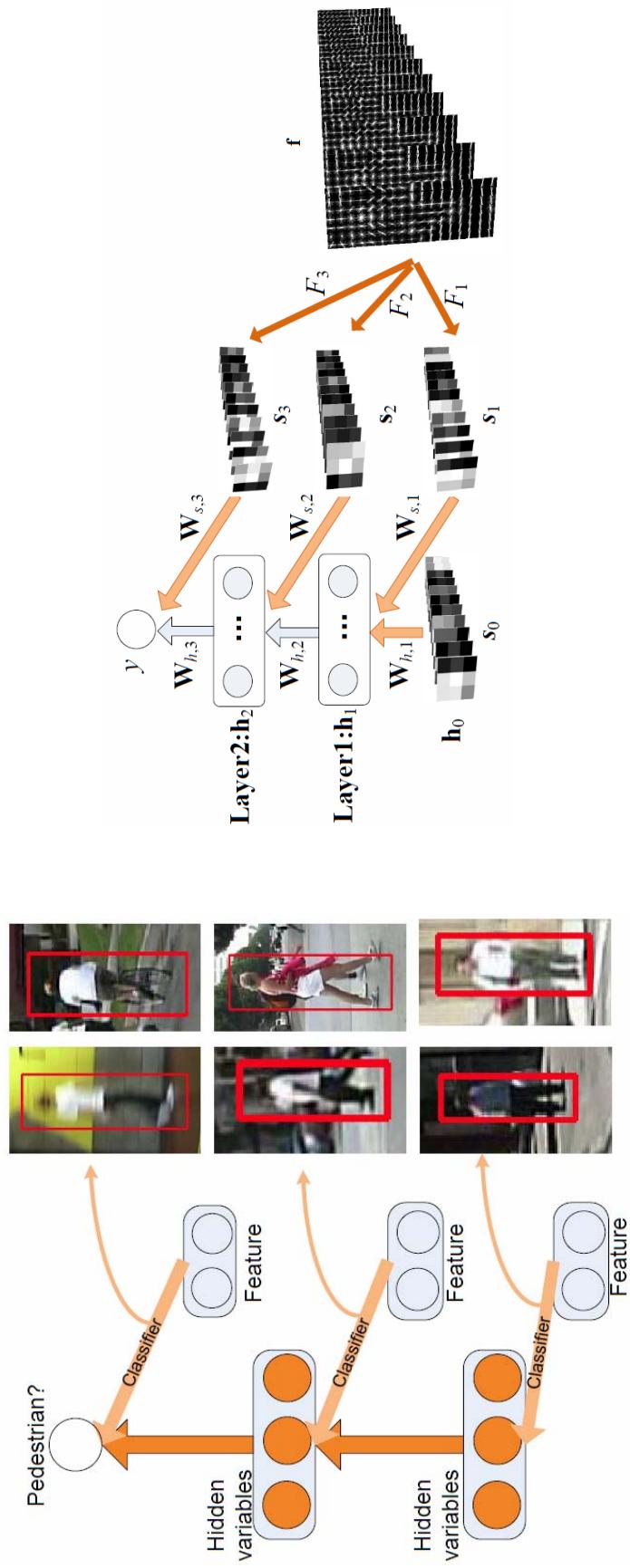
Motivated by Cascaded Classifiers and Contextual Boost

- The classifier of each stage deals with a specific set of samples
- The score map output by one classifier can serve as contextual information for the next classifier



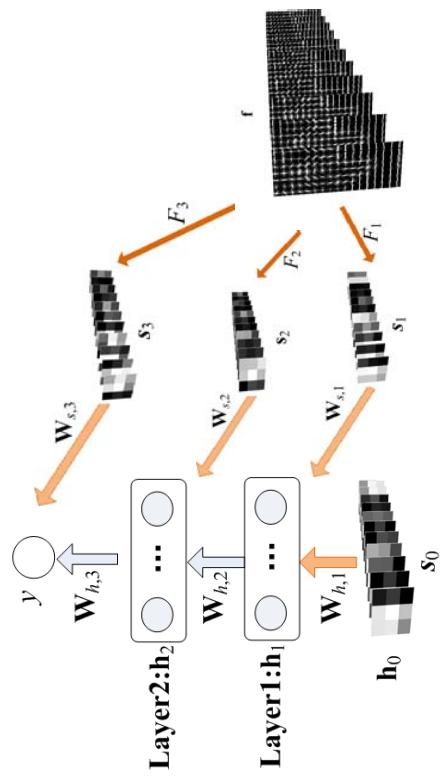
Conventional cascaded classifiers for detection

- Our deep model keeps the score map output by the current classifier and it serves as contextual information to support the decision at the next stage
- Cascaded classifiers are jointly optimized instead of being trained sequentially
- To avoid overfitting, a stage-wise pre-training scheme is proposed to regularize optimization
- Simulate the cascaded classifiers by mining hard samples to train the network stage-by-stage



Training Strategies

- Unsupervised pre-train $\mathbf{W}_{h,i+1}$ layer-by-layer, setting $\mathbf{W}_{s,i+1} = 0, \mathbf{F}_{i+1} = 0$
- Fine-tune all the $\mathbf{W}_{h,i+1}$ with supervised BP
- Train \mathbf{F}_{i+1} and $\mathbf{W}_{s,i+1}$ with BP stage-by-stage
 - A correctly classified sampled at the previous stage does not influence the update of parameters
 - Stage-by-stage training can be considered as adding regularization constraints to parameters, i.e. some parameters are constrained to be zeros in the early training stages



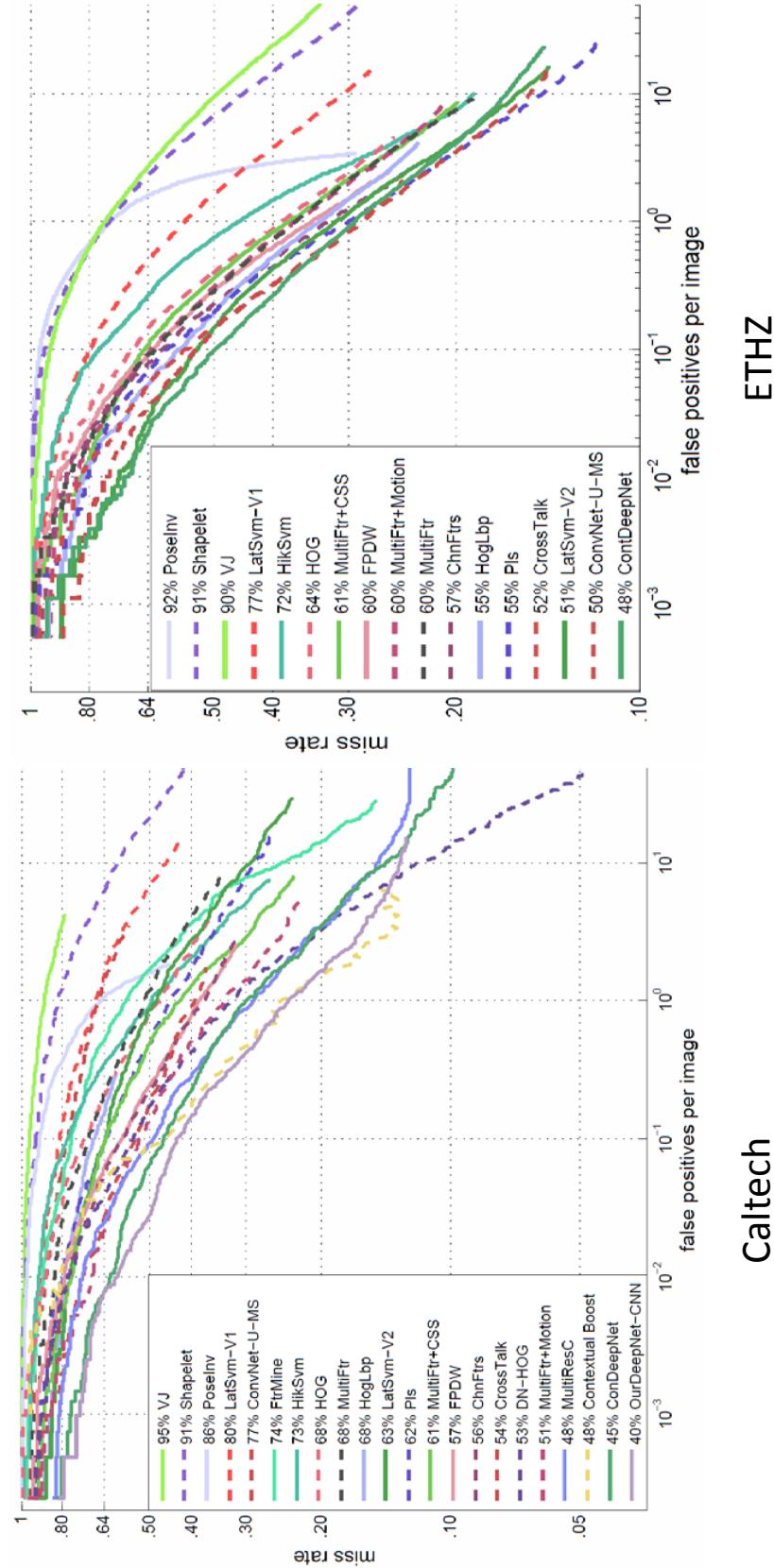
Log error function:

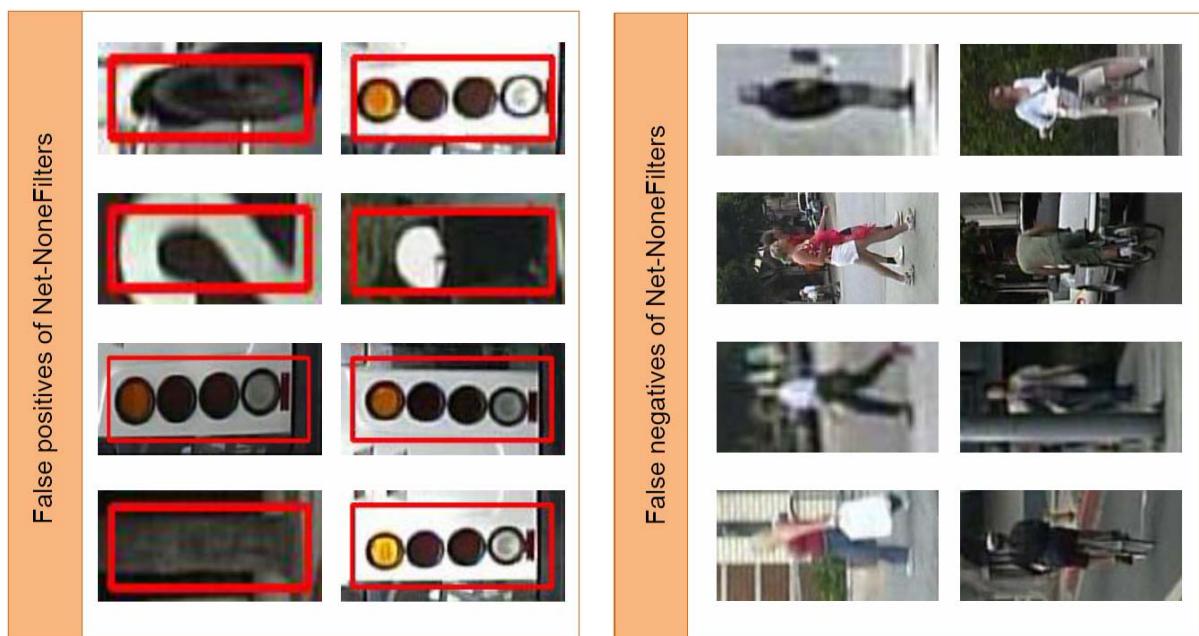
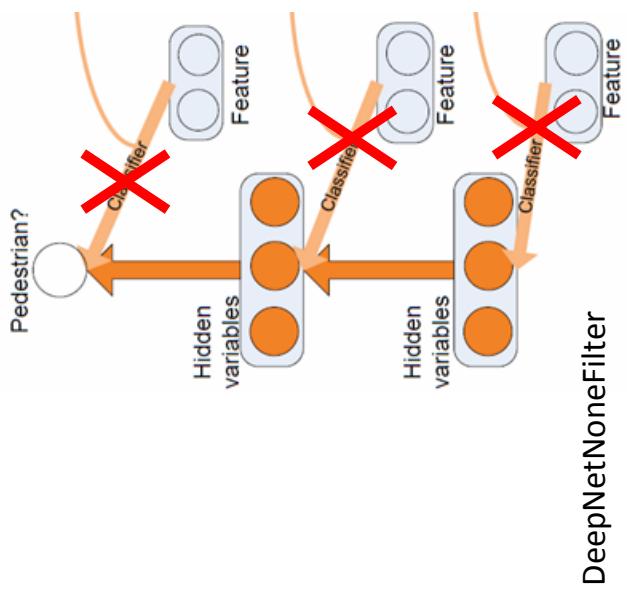
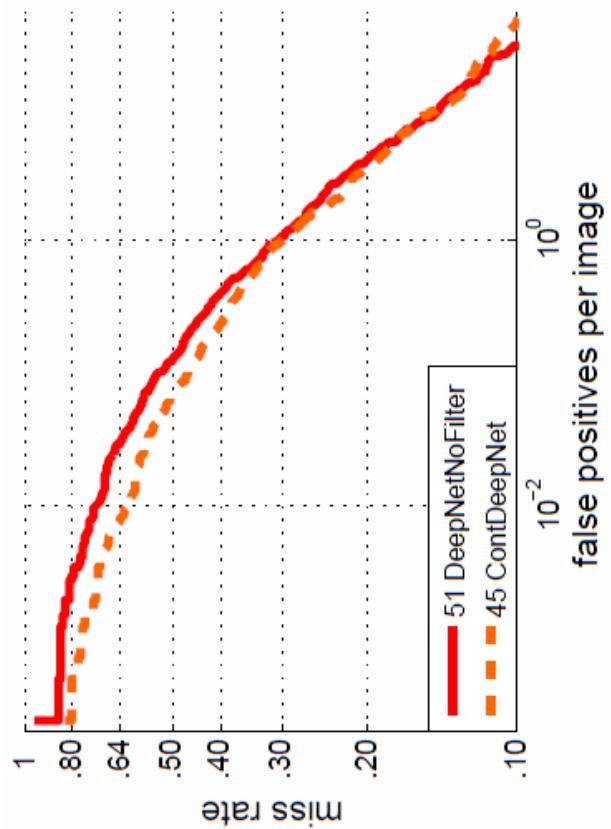
$$E = -l \log y - (1-l) \log (1-y)$$

Gradients for updating parameters:

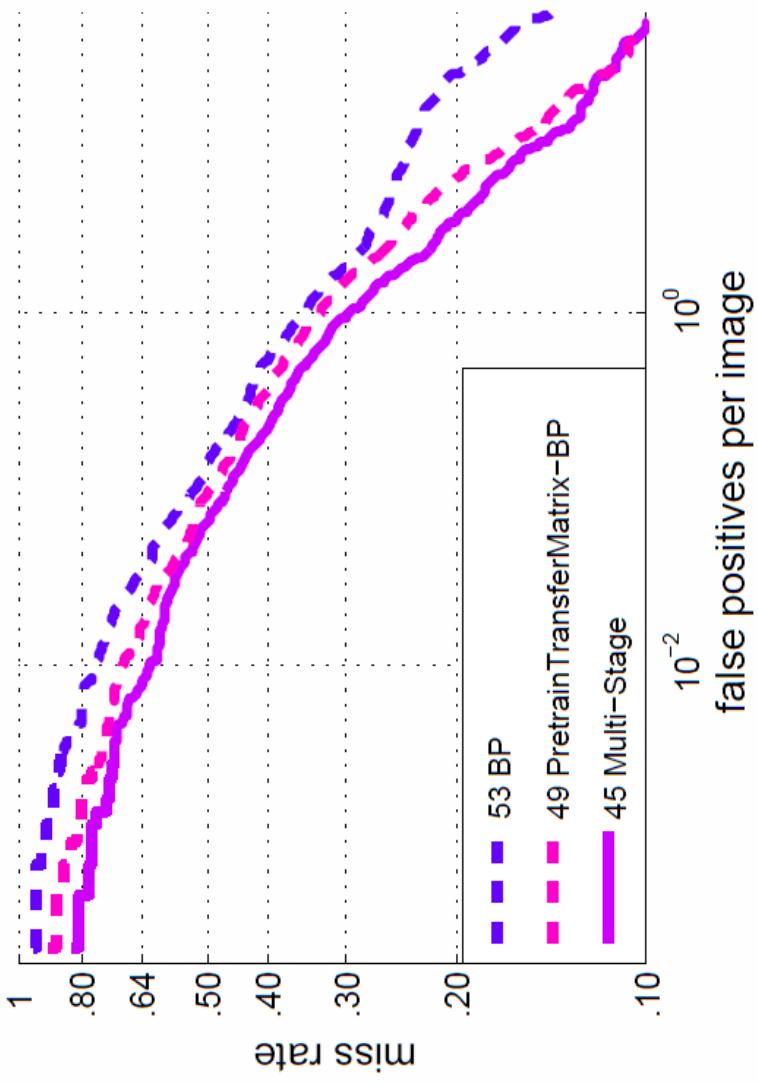
$$\frac{\partial \theta_{i,j}}{\partial E} = -\frac{\partial E}{\partial \theta_{i,j}} = -\frac{\partial E}{\partial y} \frac{\partial y}{\partial \theta_{i,j}} = -(y - l) \frac{\partial y}{\partial \theta_{i,j}}$$

Experimental Results





Comparison of Different Training Strategies

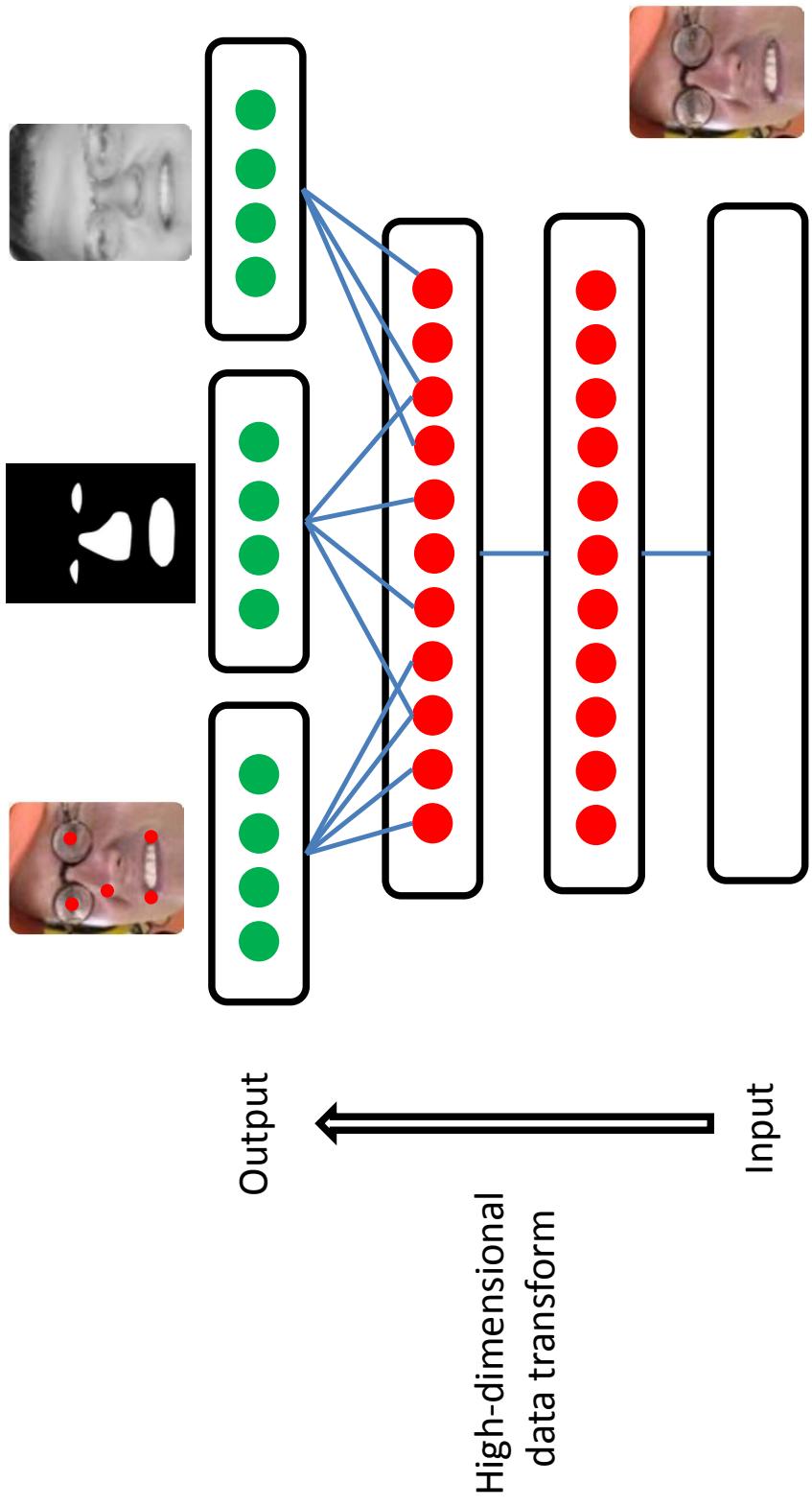


Network-BP: use back propagation to update all the parameters without pre-training

PretrainTransferMatrix-BP: the transfer matrices are unsupervised pretrained, and then all the parameters are fine-tuned

Multi-stage: our multi-stage training strategy

High-Dimensional Data Transforms



Facial keypoint detection: face image -> facial keypoint

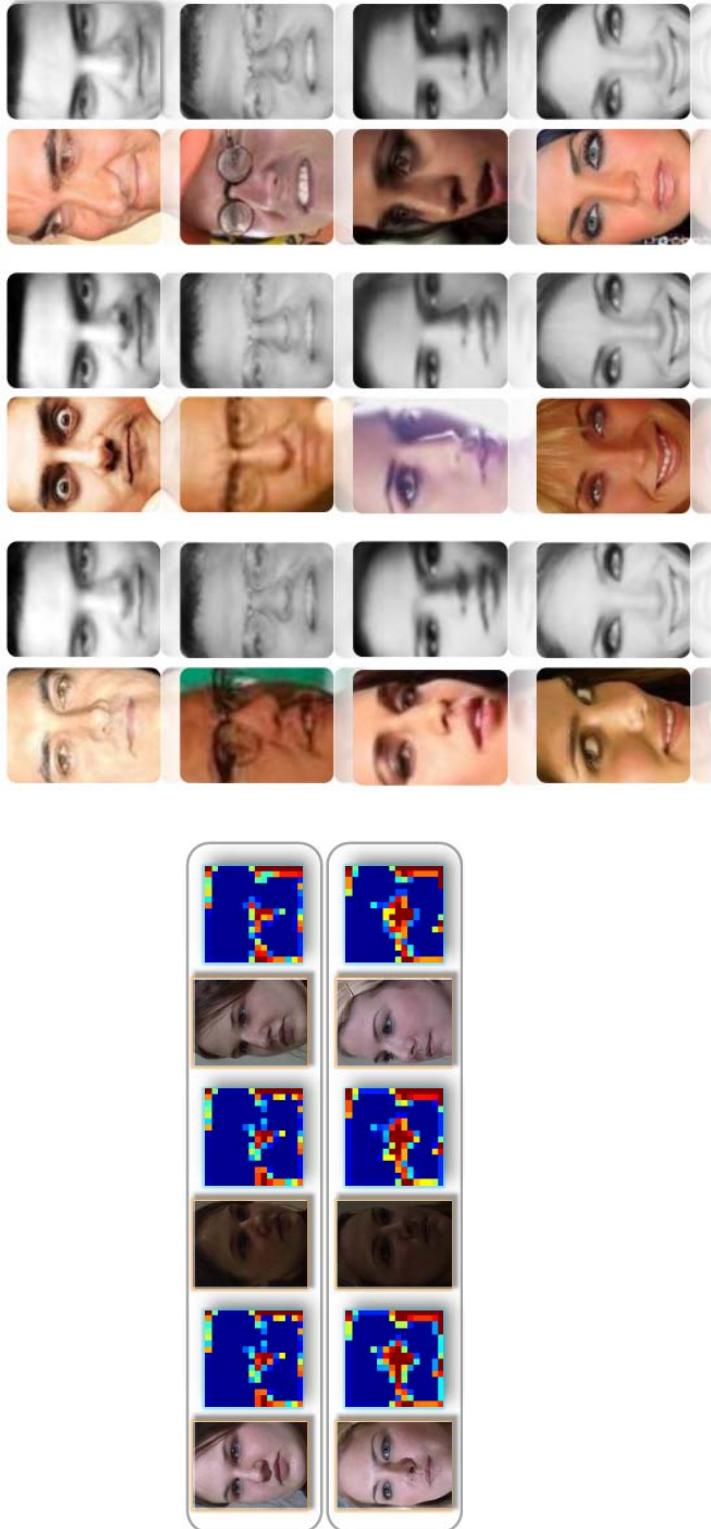
Face transform: face image in a arbitrary view -> face image in a canonical view

Face parsing: face image -> segmentation maps

Pedestrian parsing : pedestrian image -> segmentation maps

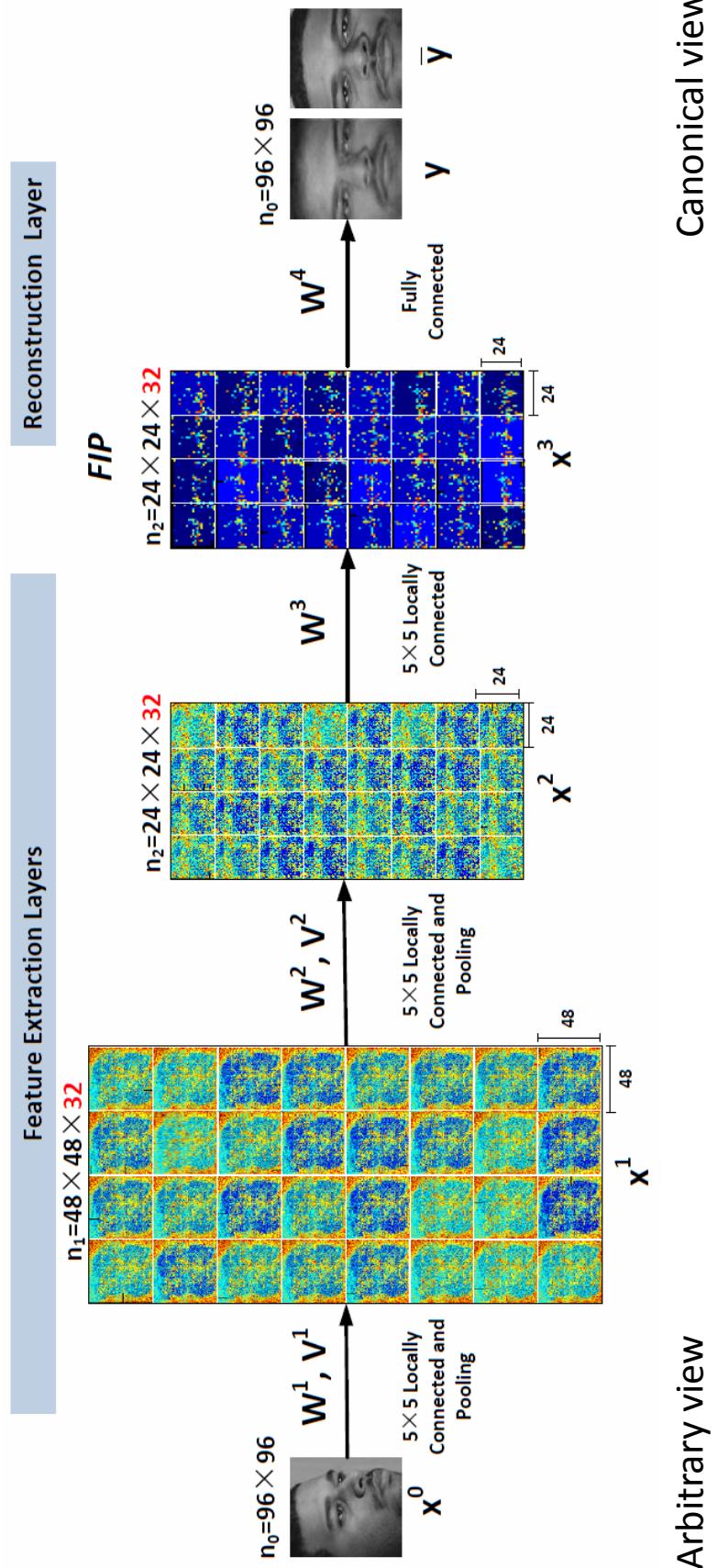
Recovering Canonical-View Face Images

- Z. Zhu, P. Luo, X. Wang, and X. Tang, “Deep Learning Identity-Preserving Face Space,” ICCV 2013.



Reconstruction examples from LFW

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





Comparison on Multi-PIE

	-45°	-30°	-15°	+15°	+30°	+45°	Avg	Pose
LGBP [26]	37.7	62.5	77	83	59.2	36.1	59.3	V
VAAM [17]	74.1	91	95.7	95.7	89.5	74.8	86.9	V
FA-EGFC[3]	84.7	95	99.3	99	92.9	85.2	92.7	X
SA-EGFC[3]	93	98.7	99.7	99.7	98.3	93.6	97.2	V
LE[4] + LDA	86.9	95.5	99.9	99.7	95.5	81.8	93.2	X
CRBM[9] + LDA	80.3	90.5	94.9	96.4	88.3	89.8	87.6	X
Ours	95.6	98.5	100.0	99.3	98.5	97.8	98.3	X

- [3] A. Asthana, T. K. Marks, M. J. Jones, K. H. Tieu, and M. Rohith. Fully automatic pose-invariant face recognition via 3d pose normalization. In *ICCV*, pages 937–944, 2011. [1](#), [5](#), [6](#)
- [4] Z. Cao, Q. Yin, X. Tang, and J. Sun. Face recognition with learning-based descriptor. In *CVPR*, pages 2707–2714, 2010. [2](#), [3](#), [6](#)
- [5] G. B. Huang, H. Lee, and E. Learned-Miller. Learning hierarchical representations for face verification with convolutional deep belief networks. In *CVPR*, pages 2518–2525, 2012. [3](#), [6](#)
- [6] W. Zhang, S. Shan, W. Gao, X. Chen, and H. Zhang. Local gabor binary pattern histogram sequence (lgphs): A novel non-statistical model for face representation and recognition. In *IICCV*, volume 1, pages 786–791, 2005. [5](#), [6](#)
- [7] S. Li, X. Liu, X. Chai, H. Zhang, S. Lao, and S. Shan. Morphable displacement field based image matching for face recognition across pose. In *ECCV*, pages 102–115, 2012. [1](#), [2](#), [5](#), [6](#)

Comparison on LFW (without outside training data)

Methods	Accuracy (%)
PLDA (Li, TPAMI'12)	90.07
Joint Bayesian (Chen, ECCV'12, 5-point align)	90.9
Fisher Vector Faces (Barkan, ICCV'13)	93.30
High-dim LBP (Chen, CVPR'13, 27-point align)	93.18
Ours (5-point align)	94.38

Comparison on LFW (with outside training data)

Methods	Accuracy (%)
Associate-Predict (Yin CVPR'12)	90.57
Joint Bayesian (Chen, ECCV'12, 5-point align)	92.4
Tom-vs-Peter (Berg, BMVC'12, 90-point align)	93.30
High-dim LBP (Chen, CVPR'13, 27-point align)	95.17
Transfer learning joint Bayesian (Cao, ICCV'13, 27-point align)	96.33
Ours (5-point align)	96.45

Face Parsing

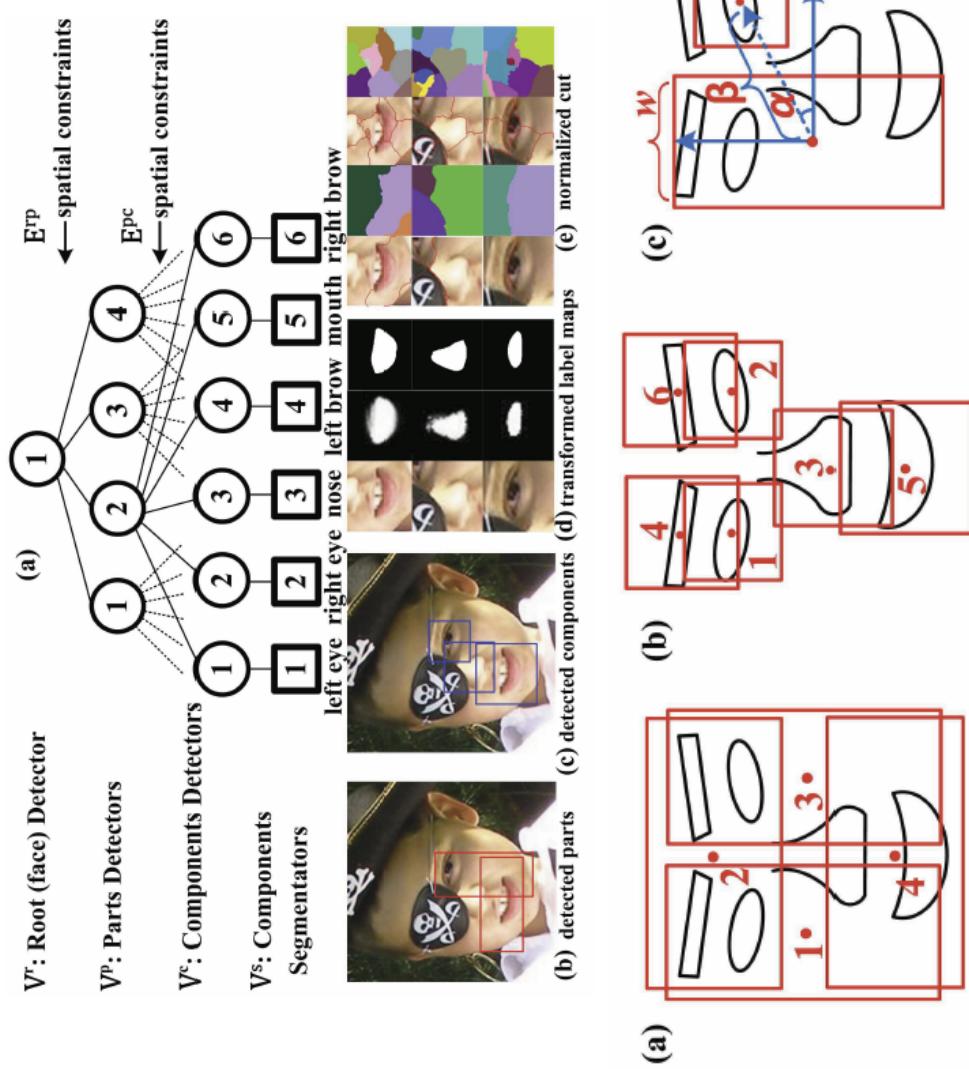
- P. Luo, X. Wang and X. Tang, “Hierarchical Face Parsing via Deep Learning,” CVPR 2012



Motivations

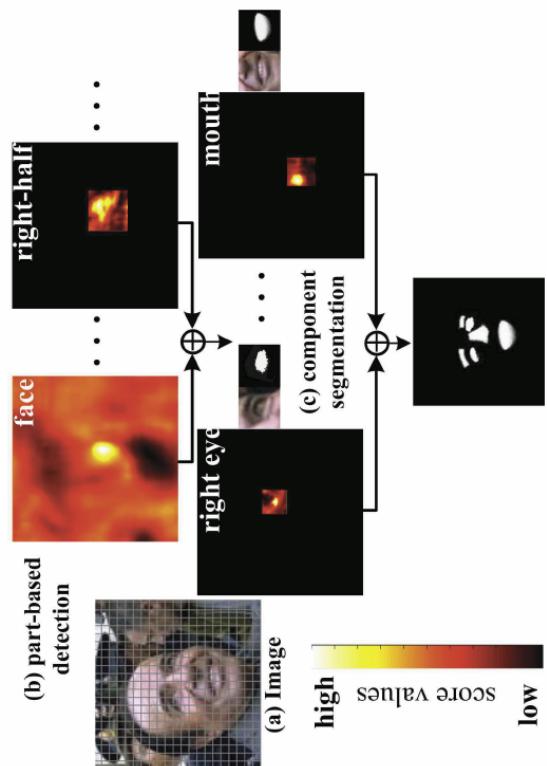
- Recast face segmentation as a cross-modality data transformation problem
- Cross modality autoencoder
- Data of two different modalities share the same representations in the deep model
- Deep models can be used to learn shape priors for segmentation

Hierarchical Representation of Face Parsing

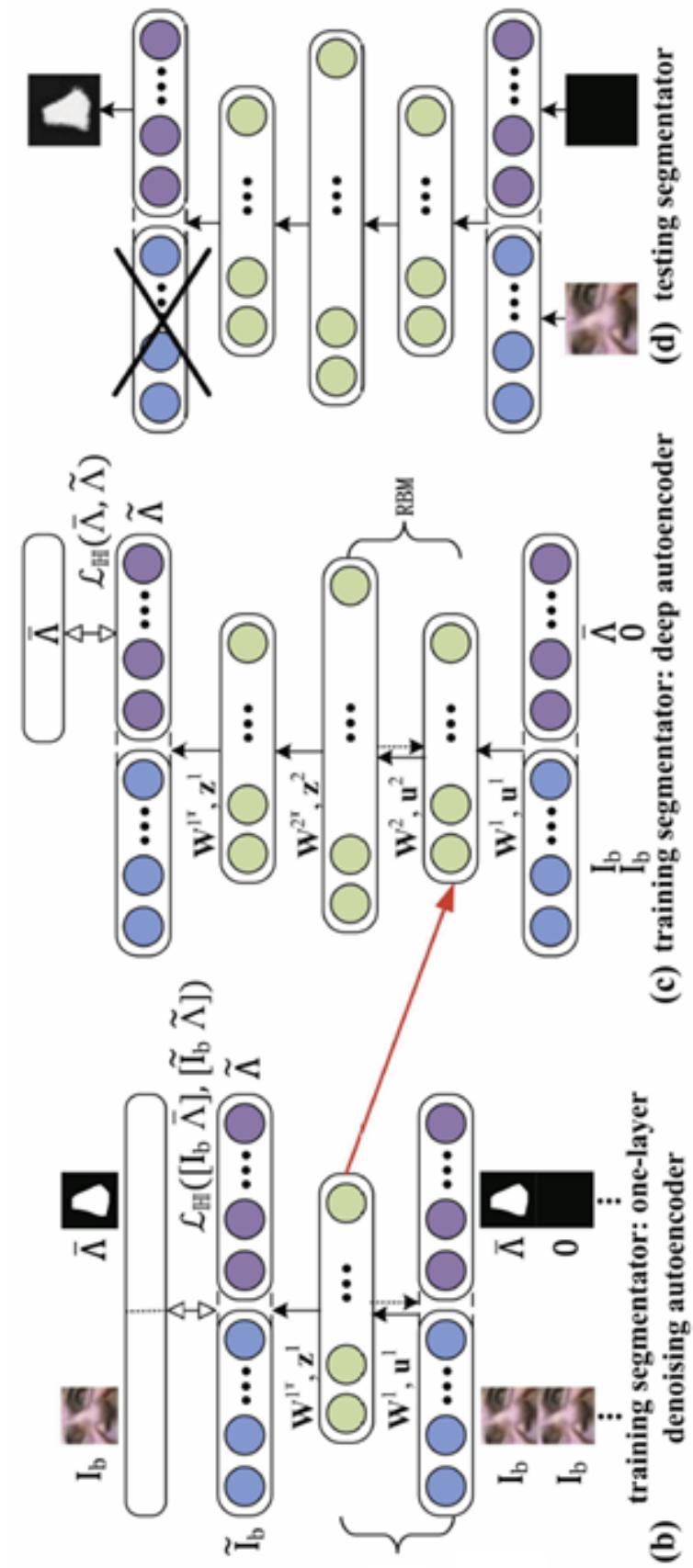


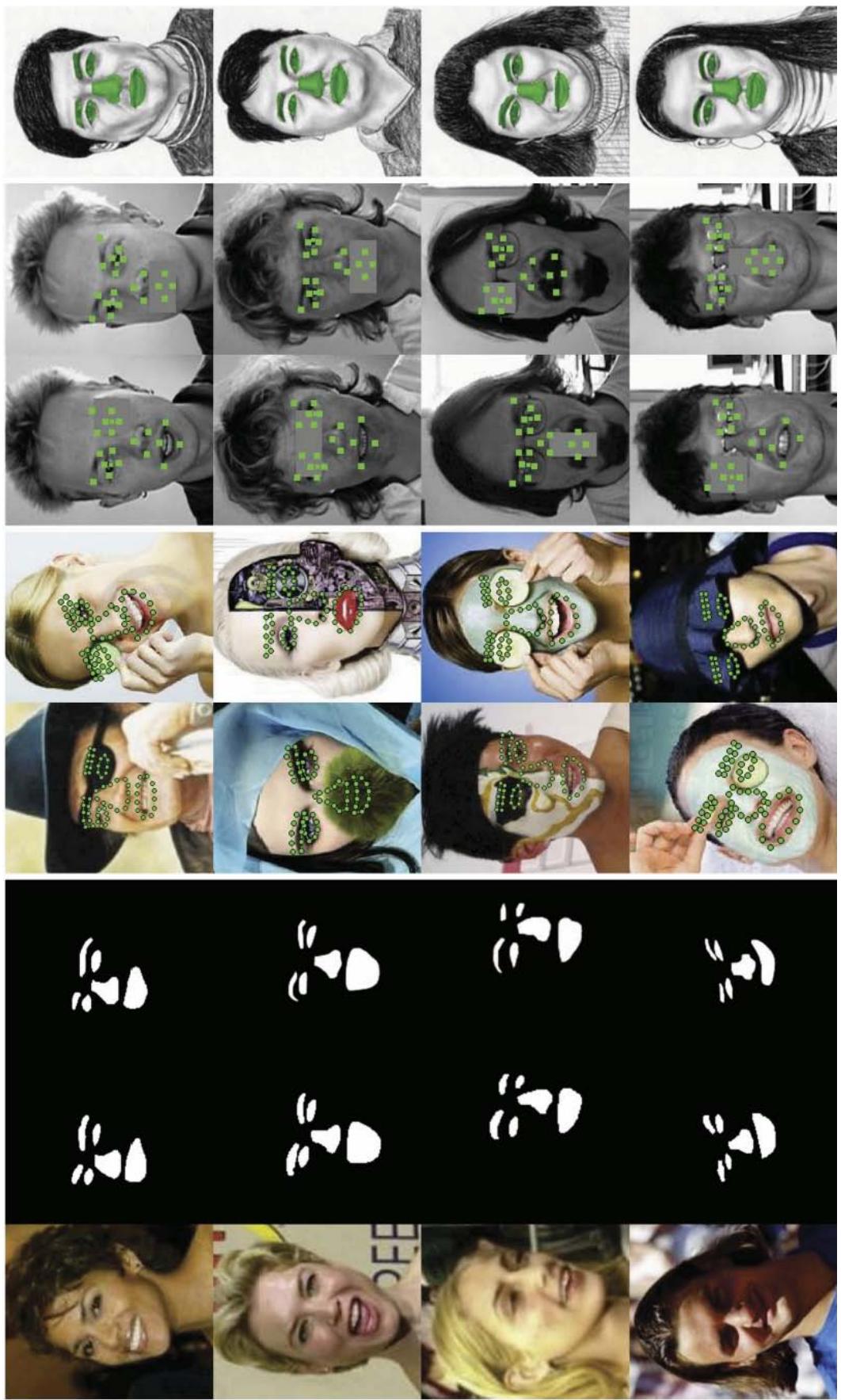
Joint Bayesian Formulation

- Detectors are trained with deep belief net (DBN) and segmentors are trained with deep autoencoder. Both have are generative models.
- Joint Bayesian framework for face detection, part detection, component detection, and component segmentation



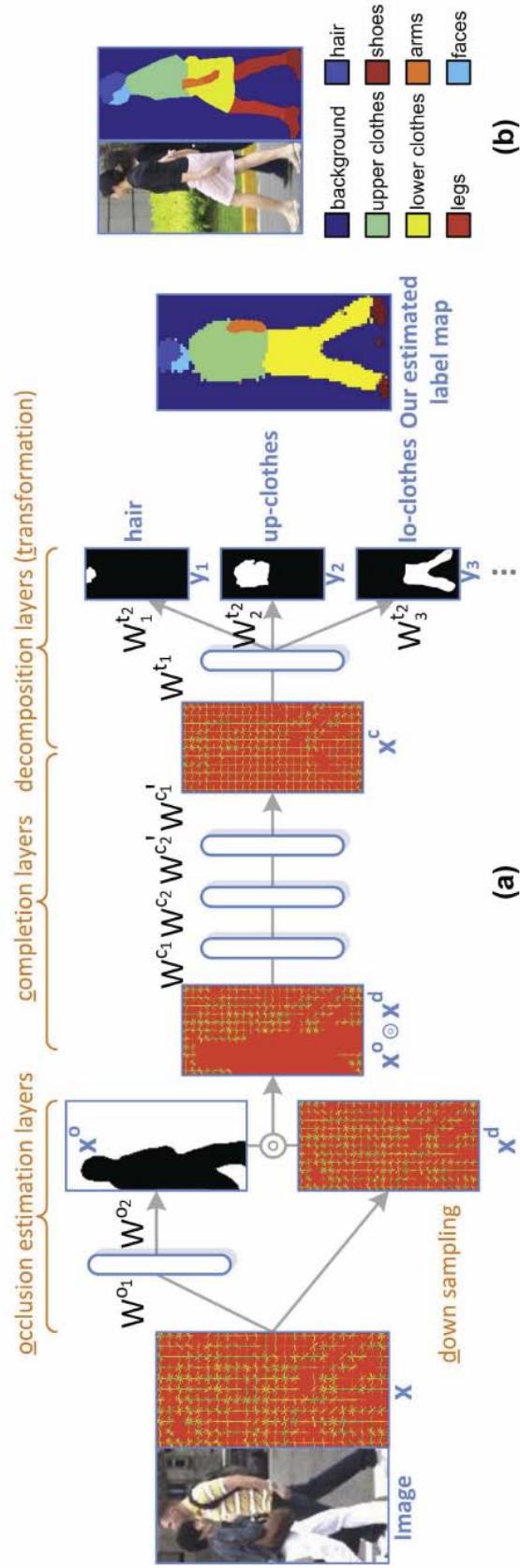
Training Segmentors

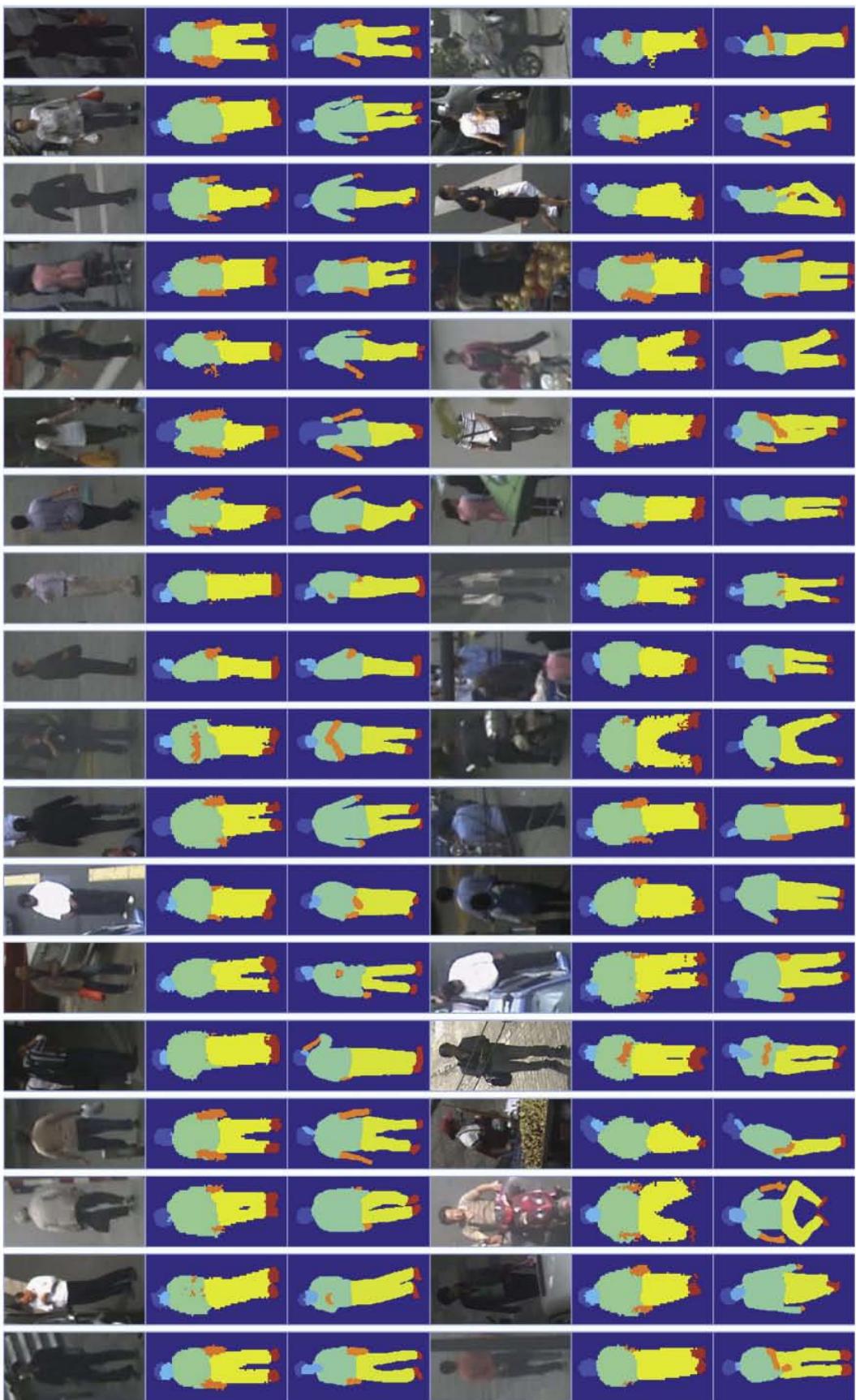




Human Parsing

- P. Luo, X. Wang, and X. Tang, “Pedestrian Parsing via Deep Decompositional Network,” ICCV 2013



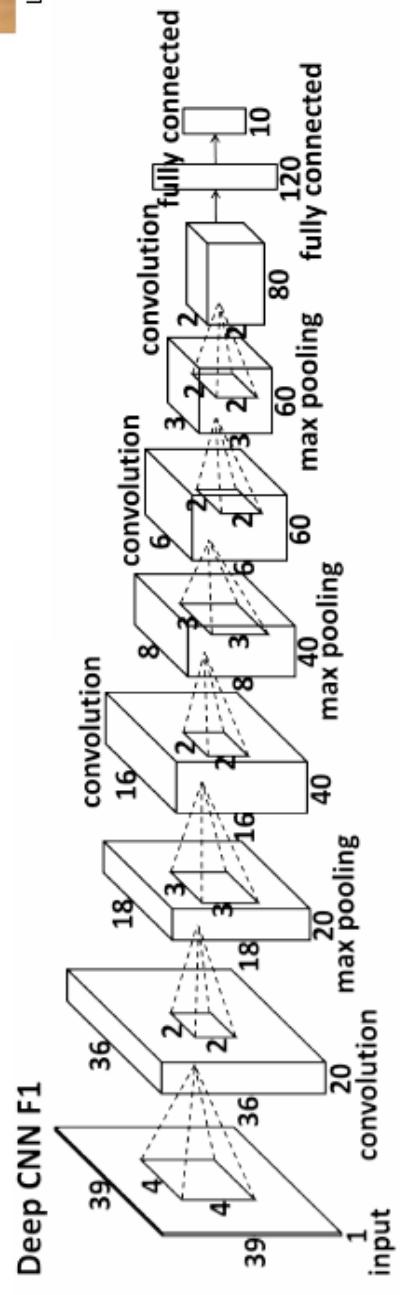
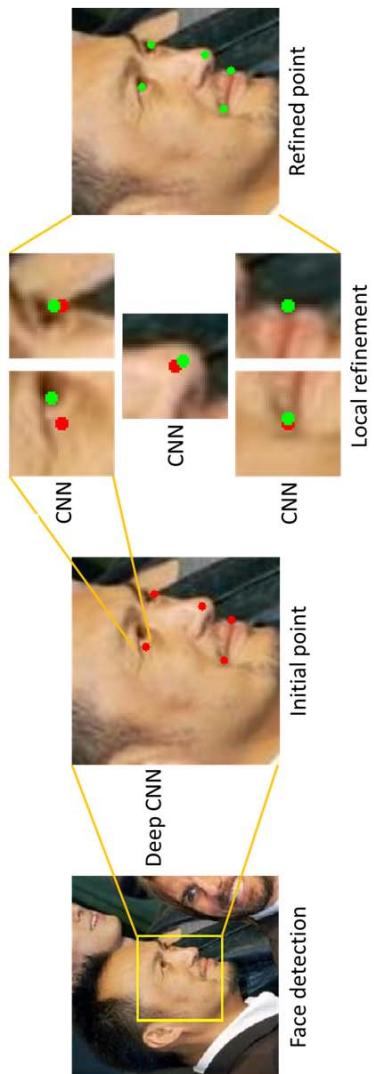
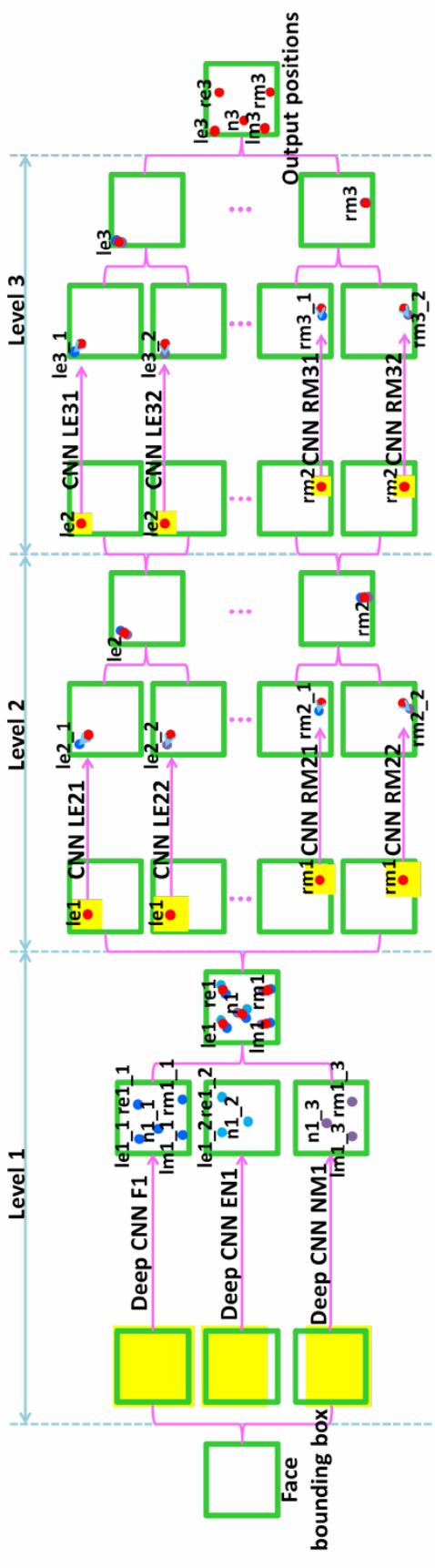


Second row: our result
Third row: ground truth

Facial Keypoint Detection

- Y. Sun, X. Wang and X. Tang, “Deep Convolutional Network Cascade for Facial Point Detection,” CVPR 2013



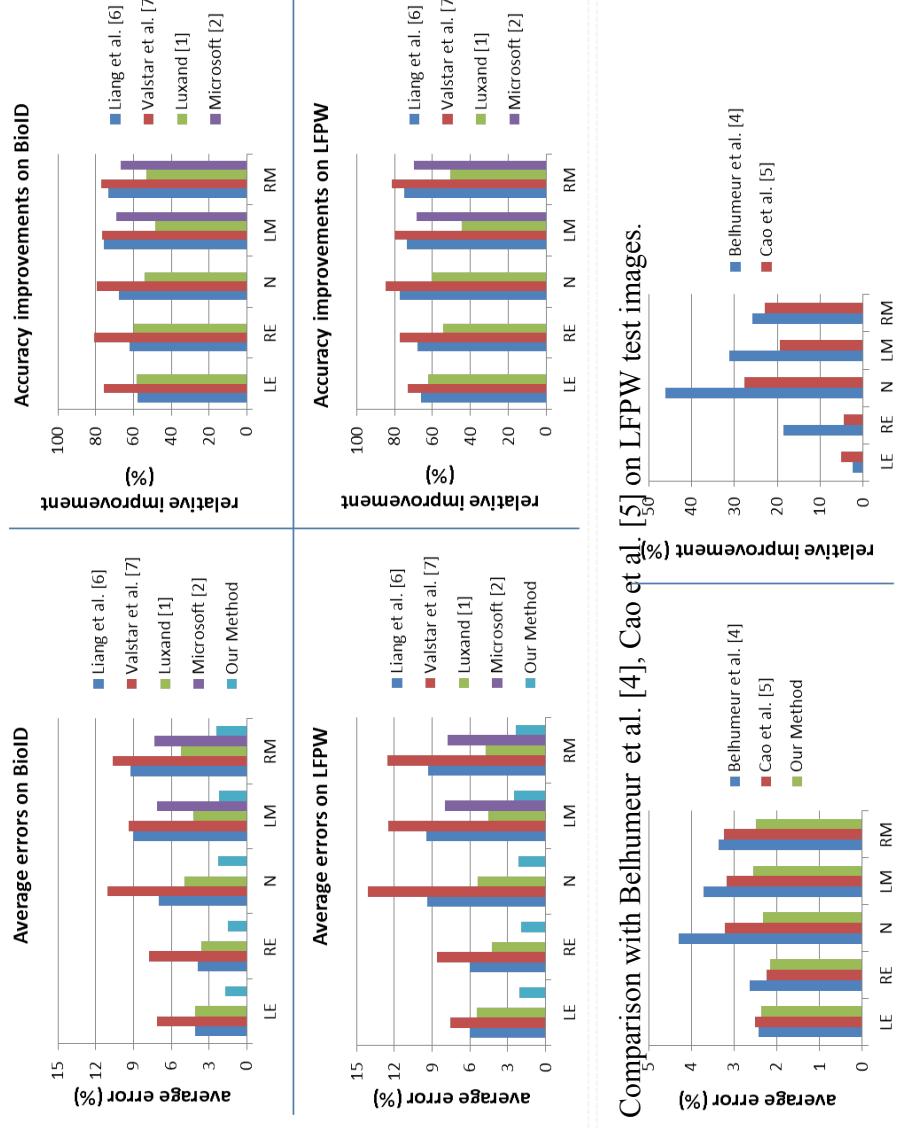


Benefits of Using Deep Model

- Take the full face as input to make full use of texture context information over the entire face to locate each keypoint
- The first network of tackling the whole face as input needs **deep structures** to extract **high-level** features
- Since the networks are trained to predict all the keypoints simultaneously, the geometric constraints among keypoints are implicitly encoded

Comparison with Liang et al. [6], Valstar et al. [7], Luxand Face SDK [1] and Microsoft Research Face SDK [2] on BioID and LFPW.

$$\text{Relative improvement} = \frac{\text{reduced average error}}{\text{average error of the method in comparison}}.$$



1. <http://www.luxand.com/facesdk/>

2. <http://research.microsoft.com/en-us/projects/facesdk/>.

3. O. Jersorsky, K. J. Kirchberg, and R. Frischholz. Robust face detection using the hausdorff distance. In Proc. AVBPA, 2001.

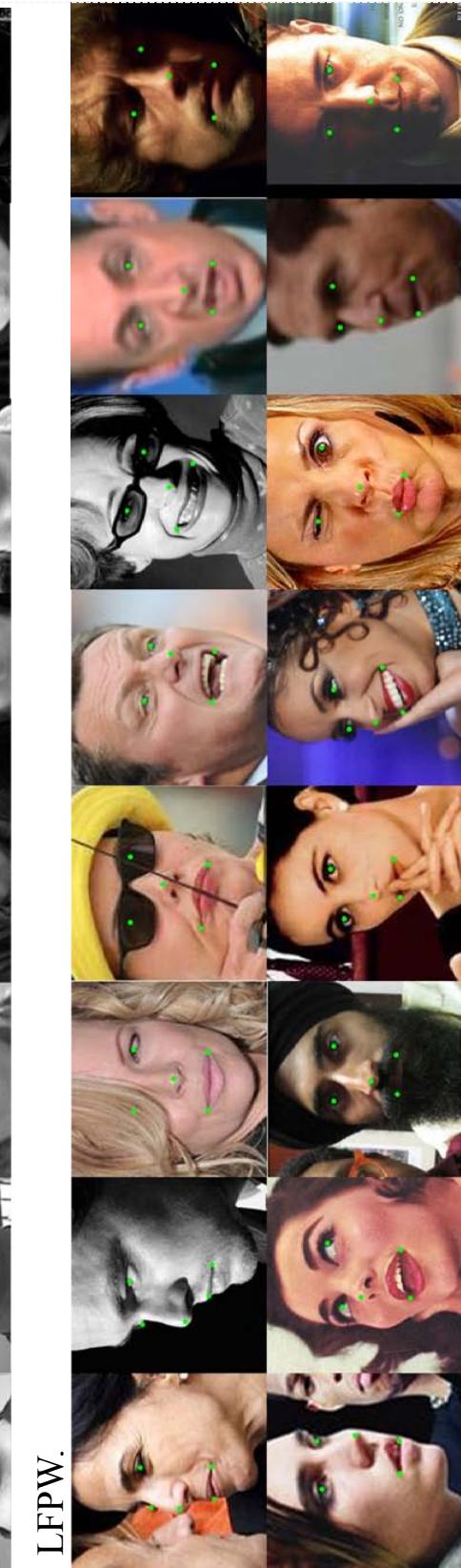
4. P. N. Belhumeur, D. W. Jacobs, D. J. Kriegman, and N. Kumar. Localizing parts of faces using a consensus of exemplars. In Proc. CVPR, 2011.

5. X. Cao, Y. Wei, F. Wen, and J. Sun. Face alignment by explicit shape regression. In Proc. CVPR, 2012.

6. L. Liang, R. Xiao, F. Wen, and J. Sun. Face alignment via component-based discriminative search. In Proc. ECCV, 2008.

7. M. Valstar, B. Martínez, X. Biñéfa, and M. Pantic. Facial point detection using boosted regression and graph models. In Proc. CVPR, 2010.

Validation.



BioID.



LFPW.



Conclusions

- Deep learning can jointly optimize key components in vision systems
- Prior knowledge from vision research is valuable for developing deep models and training strategies
- Deep learning can solve some vision challenges as problems of high-dimensional data transform
- Challenging prediction tasks can make better use the large learning capacity and avoid overfitting

People working on deep learning in our group



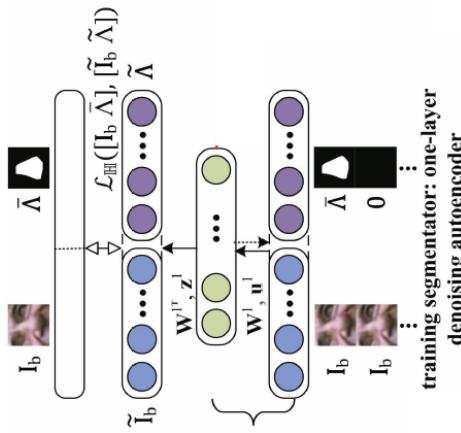
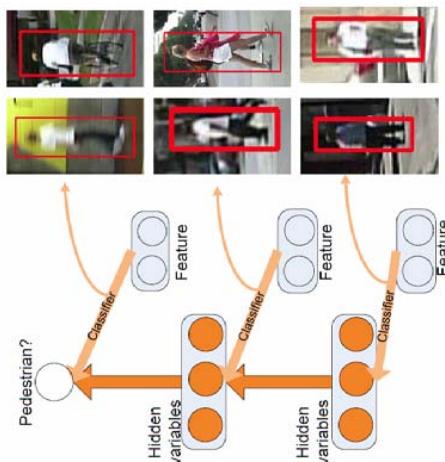
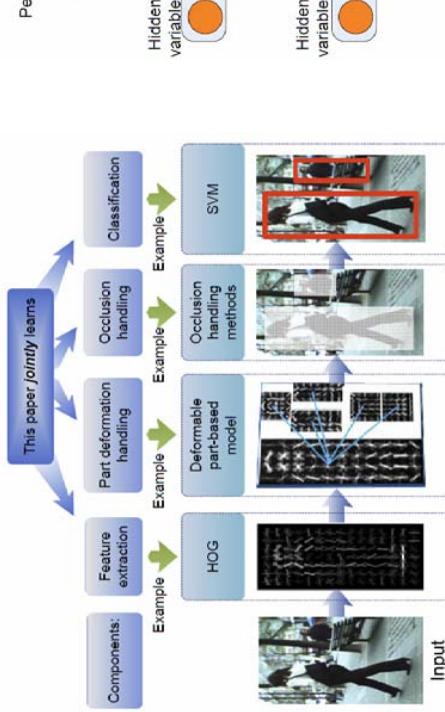
Wanli Ouyang Ping Luo Yi Sun Xingyu Zeng Zhenyao Zhu

Acknowledgement

Hong Kong Research Grants Council

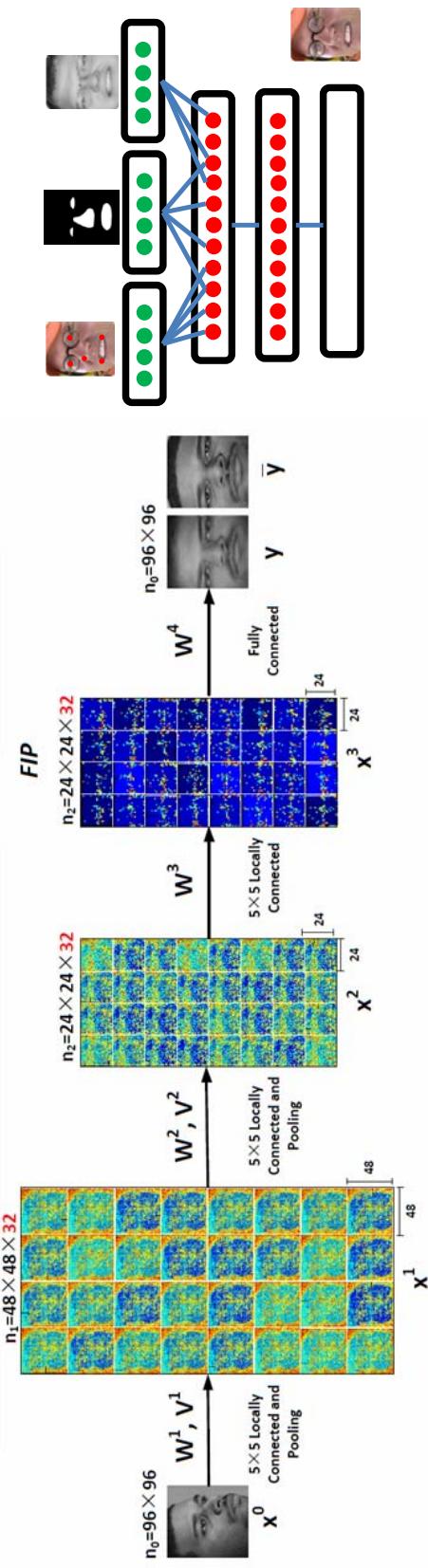
中国自然科学基金

Thank you!



Reconstruction Layer

Feature Extraction Layers



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

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